

Stable Generalized Finite Element Method and associated iterative schemes; application to interface problems

Kenan Kergrene ¹Ivo Babuška ²Uday Banerjee ³

Abstract

The Generalized Finite Element Method (GFEM) is an extension of the Finite Element Method (FEM), where the standard finite element space is augmented with a space of non-polynomial functions, called the enrichment space. The functions in the enrichment space mimic the local behavior of the unknown solution of the underlying variational problem. GFEM has been successfully applied to a wide range of problems. However, it often suffers from bad conditioning, i.e., its conditioning may not be robust with respect to the mesh and in fact, the conditioning could be much worse than that of the standard FEM. In this paper, we present a numerical study that shows that if the “angle” between the finite element space and the enrichment space is bounded away from 0, uniformly with respect to the mesh, then the GFEM is stable, i.e., the conditioning of GFEM is not worse than that of the standard FEM. A GFEM with this property is called a Stable GFEM (SGFEM). The last part of the paper is devoted to the derivation of a robust iterative solver exploiting this angle condition. It is shown that the required “wall-clock” time is greatly reduced compared to popular GFEMs used in the literature.

Keywords: Generalized Finite Element Method (GFEM), Partition of Unity Method (PUM), Stable GFEM (SGFEM), Condition Number, Angle Condition

1 Introduction

The Generalized Finite Element Method (GFEM) has sparked a lot of interest in the last 20 years and has been successfully applied to a wide range of engineering problems, e.g., crack-propagation, material modeling, and solid–fluid interactions. We refer to the review articles [1, 10, 19, 21] and the citations therein for various applications of GFEM. The method has been incorporated into commercial codes, e.g., Abaqus and LS-DYNA [13, 29]. It is also known in the literature as the Extended Finite Element Method (XFEM). We will simply refer to the method as GFEM and we will address special instances of this method such as SGFEM and M-GFEM.

As the name GFEM/XFEM suggests, the GFEM is a generalization/extension of the standard Finite Element Method (FEM). Specific non-polynomial local basis functions that mimic special

¹Department of Mathematics and Industrial Engineering, École Polytechnique de Montréal, Canada.

²ICES, University of Texas at Austin, Austin, TX, United States.

³Department of Mathematics, 215 Carnegie, Syracuse University, Syracuse, NY 13244, United States. E-mail: banerjee@syr.edu.

features (e.g., singularity) of the solution of the underlying PDE model of interest, are used in this method in addition to the standard “hat-functions.” These additional local basis functions are called the local *enrichment* functions. In fact, the GFEM is a particular instance of the Partition of Unity Method (PUM). The PUM, developed in [6, 31, 32], allows the use of any Partition of Unity (PU) together with local enrichment functions. The GFEM is a PUM, where the finite element “hat-functions” serve as the PU. Various methods for solving multi-scale problems are also based on PUM; see for example [4, 5, 19]. The PUM based on a “flat-top” PU was developed in [22, 38]. A similar idea, referred to as $h - p$ Cloud method, was developed in [17, 18]. The original idea of GFEM, i.e., the use of hat-functions as the PU, was introduced in [3]. Since then, the GFEM has been developed, refined, and used in various applications in two and three dimensions, e.g., in [9, 14, 15, 16, 36, 41, 42, 43, 44]. The GFEM is often referred to as the XFEM in the literature. We mention that the use of local non-polynomial approximation, not in the framework of PUM, was suggested earlier in [7].

The XFEM/GFEM was initially developed as a computational method with essentially intuitive understanding of the necessity of appropriate enrichments for convergence. The method was primarily tested numerically to ensure convergence. Appropriate enrichment functions for various applications were identified in the literature and the optimal convergence of the approximate solution was shown through computations. A rigorous mathematical proof of optimal convergence was derived in [35], in the context of a crack problem.

Though approximability and optimal convergence are very important features of a numerical method such as GFEM, it is equally important that the underlying linear system could be solved accurately and efficiently. Solving such linear systems accurately and efficiently depends on the stability of the GFEM, i.e., on the conditioning of the underlying linear system. It was reported early in [6, 21] that the GFEM could be unstable and that its conditioning may not be robust with respect to the mesh. However, there are very few papers that addressed these issues by carefully studying the conditioning of the GFEM and by examining the performance of associated iterative solvers to solve the linear system. Various ad-hoc stabilization procedures were used to address these issues in [8, 30, 33, 41, 43]. Stabilization based on a local orthogonalization idea was used in [30]. We mention however that local orthogonalization was also addressed in the context of PUM with flat-top PU in [39].

Extensive literature is available on the loss of accuracy in the computed solution of a linear system; we refer to the monographs [26, 45]. In [2], it was shown that the Scaled Condition Number (SCN) of the matrix, related to FEM and GFEM, is a good indicator of the stability and loss of accuracy in the solution obtained from elimination methods.

The conditioning of GFEM was addressed in [2] where the idea of a stable GFEM (SGFEM) was introduced. In general, a GFEM is called stable (SGFEM) if

- (i) it yields the optimal order of convergence, and

- (ii) the SCN of the linear system associated with the GFEM is of the same order $O(h^{-2})$ (h being the discretization parameter) as that of a standard FEM in a robust manner with respect to the mesh.

It was shown in [2] that the SCN of a GFEM could be much higher than that of the FEM, e.g., $O(h^{-4})$. It was mathematically established in that paper that if the enrichments satisfy two specific conditions, then the SCN of the underlying GFEM is of the same order as that of a standard FEM. For various problems in the 1-D setting, a simple modification based on subtracting the piecewise linear interpolant of the standard enrichment was suggested in [2] and it was shown that the modified GFEM was indeed an SGFEM for these problems. However, the modification suggested in [2] may lead to loss of accuracy in some problems in higher dimensions as shown in [23, 24, 37]. It was shown in [23, 24] that a further modification of “Heaviside enrichment,” in the context of a problem with a crack, is required for a GFEM to be an SGFEM, i.e., the further modification restores the accuracy of the computed solution while retaining the well-conditioning of the linear system. Thus a GFEM with the simple modification of enrichments as suggested in [2] may not yield an SGFEM for every problem; further modifications of the enrichments may be required for a GFEM to be an SGFEM.

In this paper, we consider an “interface problem” modeled by a scalar second order elliptic PDE in 2-D with piecewise smooth coefficients. We will numerically investigate the accuracy, conditioning, and the robustness of the GFEM associated with various forms of enrichments used in the literature, when applied to this problem. We will especially investigate the performance of an iterative procedure to solve the underlying linear system of the GFEM, where the stopping criterion is based on computed *discretization error* and *truncation error*.

In particular, we will consider GFEM with (i) the “topological enrichment” where a minimal order of enriched nodes are used, (ii) M-GFEM which is a generalization of topological enrichment, (iii) the “geometrical enrichment” where all the nodes in a fixed neighborhood of the interface are enriched, and (iv) the so-called SGFEM obtained by the simple modification of M-GFEM enrichments, as suggested in [2]. Through numerical experiments, we show that the SGFEM is accurate for interface problems and that it does not lose accuracy as suggested in [23, 24, 37]; thus no further modification of the enrichment is required for the interface problem (in contrast with [23, 24] where it was required to restore accuracy). We also show that among all the enrichments considered in the paper, the SGFEM is the only method that is well-conditioned and robust with respect to the mesh for the interface problem. These properties of the SGFEM will be proved mathematically in a forthcoming paper. One of the most important features of the current paper is the study of the performance of an iterative procedure to solve the linear system associated with M-GFEM and SGFEM. For a given error tolerance τ , we have computed the solutions of the linear systems for the M-GFEM and SGFEM using the iterative method with the stopping criterion based jointly on discretization–truncation errors, as mentioned before. We observed that SGFEM requires

fewer iterations and less “wall-clock” time than the M-GFEM.

The outline of the paper is as follows: in Section 2 we define the interface problem. In Section 3 we describe the GFEMs with various enrichments together with their convergence and conditioning properties in 1-D to communicate the idea in a simpler setting. In Section 4 we consider a “straight interface” problem in 2-D with no singularity, which could be viewed as a “laboratory problem.” We describe the GFEM with various enrichments for the straight interface problem, define an “angle-condition” that dictates the conditioning of the GFEM, and present various numerical results addressing the accuracy, conditioning, as well as the relation between the angle condition and conditioning. In this section we also discuss the notion of robustness with respect to the mesh. The numerical results clearly indicate that the GFEM with modified M-GFEM enrichments is indeed an SGFEM and is the most robust of all the methods considered. In Section 5 we present similar numerical results for a circular interface problem and come to the same conclusions as in Section 4. In Section 6 we describe the iterative method and the stopping criteria. We study its performance on linear systems associated with FEM, M-GFEM, and SGFEM.

2 Formulation of the interface problem

Let $\Omega \subset \mathbb{R}^2$ be a bounded, simply connected domain with smooth boundary $\partial\Omega$. Consider another simply connected domain $\Omega_1 \subset \Omega$ and set $\Omega_0 := \Omega \setminus \overline{\Omega}_1$. $\Gamma := \overline{\Omega}_0 \cap \overline{\Omega}_1$ is called the *interface* and is assumed to be smooth.

We are interested in the weak solution u of the problem

$$-\nabla \cdot (a \nabla u) = f + q\delta(\Gamma), \quad \text{in } \Omega, \quad (2.1)$$

with the boundary condition

$$u = g_D, \quad \text{on } \partial\Omega, \quad (2.2)$$

$$\text{or,} \quad a \nabla u \cdot \vec{n} = g_N, \quad \text{on } \partial\Omega, \quad (2.3)$$

where $0 < \beta_0 \leq a(\mathbf{x}) \leq \beta_1$, $a_i(\mathbf{x}) := a|_{\Omega_i}$, $f_i(\mathbf{x}) := f|_{\Omega_i}$ are smooth functions on $\overline{\Omega}_i$ for $i = 0, 1$; in particular, $a_i \in C^2(\overline{\Omega}_i)$ and $f_i \in C(\overline{\Omega}_i)$. Moreover, $\delta(\Gamma)$ is the Dirac function on Γ , $q(s)$ is smooth on Γ , $g_D(s)$ is smooth on $\partial\Omega$ and $g_N(s)$ is smooth on $\partial\Omega \cap \overline{\Omega}_i$; in particular, $g_D(s) \in C^2(\partial\Omega)$, $g_N(s) \in C^1(\partial\Omega \cap \overline{\Omega}_i)$, and $q(s) \in C^1(\Gamma)$, where s is the arc length parameter.

The weak solution $u \in H^1(\Omega)$, $u|_{\partial\Omega} = g_D$ of the Dirichlet boundary value problem (2.1),(2.2) satisfies

$$B(u, v) := \int_{\Omega} a \nabla u \cdot \nabla v \, d\mathbf{x} = \int_{\Omega} f v \, d\mathbf{x} + \int_{\Gamma} q v \, ds, \quad \text{for all } v \in H_0^1(\Omega).$$

The solution of the above variational problem exists and is unique.

In the case of Neumann boundary value problem (2.1),(2.3), the weak solution $u \in H^1(\Omega)$ satisfies

$$B(u, v) = \int_{\Omega} f v \, d\mathbf{x} + \int_{\partial\Omega} g_N v \, ds + \int_{\Gamma} q v \, ds, \quad \text{for all } v \in H^1(\Omega). \quad (2.4)$$

The solution of the above variational problem exists and is unique up to an additive constant, provided the data f, g_N, q satisfy the compatibility condition

$$\int_{\Omega} f \, d\mathbf{x} + \int_{\partial\Omega} g_N \, ds + \int_{\Gamma} q \, ds = 0. \quad (2.5)$$

Under the assumed smoothness on the input data f, q, g_D, g_N , the solution u of the variational problem for (2.1)–(2.3) is continuous on $\overline{\Omega}$ and $u \in C^2(\overline{\Omega}_i)$ for $i = 0, 1$. In particular, the solution does not have any singularity anywhere in $\overline{\Omega}$. We mention that in this paper we do not address the minimum regularity requirements on the input data for the analysis of GFEM.

Remark 2.1 Note that the data $q(s)$ appears in the right hand side of the variational problem and has no effect on the choice of enrichments in the GFEM and on the features of GFEM that we study in this paper. Thus, we will use $q(s) = 0$ in all our numerical experiments.

In the numerical experiments presented in this paper, we will consider $\overline{\Omega}_1$ to be a closed disk, the interface Γ is thus a circle, away from the boundary $\partial\Omega$ (see Figure 9). We will also consider a different interface problem, namely, the “straight interface problem,” where the interface Γ is a straight line intersecting the boundary $\partial\Omega$ at two points (see Figure 2). In general, the solution of this problem will have singularities at the points $\Gamma \cap \partial\Omega$. However, we will consider a manufactured solution of the straight interface problem that does not have any singularities, mimicking the property of the solution of a closed interface with data described before. We chose the straight interface problem to highlight in an easy and efficient manner the robustness properties of various GFEMs that we consider in this paper. Note that the results related to the straight interface problem are directly relevant to the problems with non-circular interfaces, where part of $\partial\Omega_1$ is a straight line. Moreover, we will consider only the Neumann problem (2.4) in all our numerical experiments.

3 Various GFEMs on a 1-D problem

The goal of this paper is to study certain features of various GFEMs, e.g., the accuracy and conditioning when applied to an interface problem. We describe these methods and the associated features for a 1-D problem, which will allow us to communicate the main ideas in a simpler setting.

Let $\Omega = (0, 1)$, $\Omega_0 = (0, \gamma)$, $\Omega_1 = (\gamma, 1)$, where $0 < \gamma < 1$ is the interface. We consider the

boundary value problem

$$\begin{aligned} -(au')' &= f, \quad \text{in } \Omega, \\ u(0) &= 0, \quad au'(1) = g, \end{aligned}$$

where $f = 1$, $g = 2$ and $a(x)|_{\Omega_0} = a_0$, $a(x)|_{\Omega_1} = a_1$ with a_0, a_1 as strictly positive constants.

The weak formulation of the above problem is

$$\begin{aligned} u &\in \mathcal{E} := \{u \in H^1(\Omega) : u(0) = 0\}, \\ B(u, v) &= F(v), \quad \text{for all } v \in \mathcal{E}, \end{aligned} \tag{3.1}$$

where

$$B(u, v) := \int_{\Omega} a u' v' dx \quad \text{and} \quad F(v) := \int_{\Omega} v dx + 2v(1).$$

We denote the energy norm of $v \in \mathcal{E}$ by $\|v\|_{\mathcal{E}} := B(v, v)^{1/2}$.

The GFEM to approximate the solution of the above problem is a generalization of the standard FEM. We will describe various GFEMs below and state the results associated with their accuracy and conditioning. This section, associated with a 1-D problem, can be viewed as a conceptual synopsis of the results presented in this paper in higher dimensional interface problem.

GFEM: Let \mathcal{T}_h be the uniform mesh on Ω with nodes $x_i^h = ih$, $i \in \mathcal{N}^h := \{0, 1, \dots, m\}$ and elements $\tau_i^h = [x_{i-1}^h, x_i^h]$, $i \in \mathcal{N}_d^h := \{1, 2, \dots, m\}$, where $h = 1/m$. With each node x_i^h , $i = 1, 2, \dots, m-1$, we associate the *patch* $\omega_i^h = (x_{i-1}^h, x_{i+1}^h)$; for $i = 0, m$, we set $\omega_0^h = (x_0^h, x_1^h)$ and $\omega_m^h = (x_{m-1}^h, x_m^h)$. Clearly, $\cup_{i=0}^m \omega_i^h = \Omega$.

Let N_i^h be the usual piecewise linear “hat-functions” associated with the node x_i^h with $\text{supp}\{N_i^h\} = \overline{\omega_i^h}$ and $N_i^h(x_i^h) = 1$. Note that the interface $\gamma \in \tau_c^h$ for some c depending on h . The GFEM solution $u_h \in S^h$ satisfies

$$B(u_h, v) = F(v), \quad \text{for all } v \in S^h, \tag{3.2}$$

where the approximation space S^h is given by

$$S^h = S_{FEM}^h \oplus S_{ENR}^h = \{v = v_1 + v_2 : v_1 \in S_{FEM}^h, v_2 \in S_{ENR}^h\}, \tag{3.3}$$

where

$$\begin{aligned} S_{FEM}^h &= \text{span}\{N_i^h, i \in \mathcal{N}_d^h\}, \\ \text{and} \quad S_{ENR}^h &= \text{span}\{wN_i^h, i \in \mathcal{R}^h \subset \mathcal{N}^h\}. \end{aligned} \tag{3.4}$$

The function w in (3.4) is called the *enrichment* and S_{ENR}^h is called the *enrichment space* of the GFEM. The function w is chosen such that it mimics the exact solution. The set $\{x_i^h\}_{i \in \mathcal{R}^h}$ denotes the set of *enriched nodes*. Different GFEMs are defined by different choices of w and the choice

of \mathcal{R}^h . S_{FEM}^h is the FE space of piecewise linear functions associated with the “triangulation” \mathcal{T}_h , satisfying the homogeneous Dirichlet boundary condition at x_0^h . It is clear that $\dim(S_{FEM}^h) = m$. Note that for $w \equiv 0$, we do not consider S_{ENR}^h in the definition of S^h ; thus we have $S^h = S_{FEM}^h$ and the GFEM is just the standard FEM. GFEM is thus a generalization or extension of FEM; the approximation space S^h of GFEM is the standard FE space S_{FEM}^h augmented with the enrichment space S_{ENR}^h , as given in (3.3). In the rest of the section, we will use $x_i, \tau_i, \omega_i, N_i$ for $x_i^h, \tau_i^h, \omega_i^h, N_i^h$ with an understanding that they depend on h .

The GFEM solution $u_h \in S^h$ is obtained in the form $u_h = \sum_{i \in \mathcal{N}_d^h} c_{1,i} N_i + \sum_{k \in \mathcal{R}^h} c_{2,k} w N_k$ by solving the linear system

$$\hat{\mathbf{A}} \hat{\mathbf{c}} = \hat{\mathbf{f}}, \quad (3.5)$$

where

$$\hat{\mathbf{A}} = \begin{bmatrix} \hat{\mathbf{A}}_{11} & \hat{\mathbf{A}}_{12} \\ \hat{\mathbf{A}}_{12}^T & \hat{\mathbf{A}}_{22} \end{bmatrix}, \quad \hat{\mathbf{c}} = \begin{bmatrix} \hat{c}_1 \\ \hat{c}_2 \end{bmatrix}, \quad \hat{\mathbf{f}} = \begin{bmatrix} \hat{f}_1 \\ \hat{f}_2 \end{bmatrix},$$

and

$$\begin{aligned} \hat{\mathbf{A}}_{11} &= [B(N_j, N_i)]_{i,j \in \mathcal{N}_d^h}, \quad \hat{\mathbf{A}}_{22} = [B(w N_k, w N_l)]_{k,l \in \mathcal{R}^h}, \\ \hat{\mathbf{A}}_{12} &= [B(w N_k, N_i)]_{i \in \mathcal{N}_d^h; k \in \mathcal{R}^h}, \quad \hat{f}_1 = [F(N_i)]_{i \in \mathcal{N}_d^h}, \quad \hat{f}_2 = [F(w N_l)]_{l \in \mathcal{R}^h}, \\ \hat{c}_1 &= [c_{1,i}]_{i \in \mathcal{N}_d^h}, \quad \hat{c}_2 = [c_{2,k}]_{k \in \mathcal{R}^h}. \end{aligned}$$

Note that the matrix $\hat{\mathbf{A}}_{11}$ is the standard FEM stiffness matrix. The matrices $\hat{\mathbf{A}}_{12}$ and $\hat{\mathbf{A}}_{22}$ depend on the enrichment w . Also $\dim(\hat{\mathbf{A}}_{22}) = \text{card}(\mathcal{R}^h)$. We further note that depending on w the diagonal elements of $\hat{\mathbf{A}}_{22}$ could be small and therefore, instead of solving the linear system (3.5), we solve the linear system

$$\mathbf{A} \mathbf{x} = \mathbf{f}, \quad (3.6)$$

with

$$\mathbf{A} = \mathbf{D} \hat{\mathbf{A}} \mathbf{D} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{12}^T & \mathbf{A}_{22} \end{bmatrix}, \quad \mathbf{f} = \mathbf{D} \hat{\mathbf{f}}, \quad \mathbf{x} = \mathbf{D}^{-1} \hat{\mathbf{c}},$$

where \mathbf{D} is a diagonal matrix such that \mathbf{A} has unit diagonal elements.

We will now describe examples of GFEM for the interface problem, where the interface $\gamma \in \hat{\tau}_c = (x_{c-1}, x_c)$; note that $\hat{\tau}_c$ is the interior of the (closed) element $\tau_c = [x_{c-1}, x_c]$.

Geometric GFEM: The enrichment part S_{ENR}^h is defined with the enrichment $w = |x - \gamma|$ and $\mathcal{R}^h = \{i \in \mathcal{N}^h : |x_i - \gamma| \leq R\}$, where $R > 0$ is fixed and independent of h . Then the $\text{card}\{x_i\}_{i \in \mathcal{R}^h} = O(h^{-1})$, and consequently, the dimension of $\hat{\mathbf{A}}_{22}$ is $O(h^{-1})$. The idea of using \mathcal{R}^h with R fixed and independent of h was used in [8, 35, 44]; for other references see [10, 21].

Topological GFEM: The enrichment part S_{ENR}^h is defined with the enrichment $w = |x - \gamma|$ as before, however, $\mathcal{R}^h = \{x_{c-1}, x_c\}$. Note that \mathcal{R}^h is the union of nodes of the element τ_c containing the interface γ . Consequently, $\hat{\mathbf{A}}_{22}$ is a 2×2 matrix and the associated stiffness matrix is smaller

than that of the Geometric GFEM. This idea was first used in [9] in the context of crack propagation where the nodes close to the crack-tip were associated with \mathcal{R}^h . For other references see [10, 21].

M-GFEM: Let $w^* = |x - \gamma|$. We consider S_{ENR}^h with the enrichment function w that is continuous in Ω , $w = w^*$ in the element $\tau_c = [x_{c-1}, x_c]$ containing γ . It is *linear in every element other than τ_c* , and $w(x_i) = 0$ for all the nodes x_i *except* x_{c-1} and x_c . The graph of w is of the form “M”, as shown in Figure 1, where the parameters are $\gamma = [2 + 1/\pi]h$ and $h = 1/5$. We consider $\mathcal{R}^h = \{x_{c-2}, x_{c-1}, x_c, x_{c+1}\}$, which is the union of all the nodes in $\overline{\omega}_{c-1} \cup \overline{\omega}_c$; $\overline{\omega}_{c-1}, \overline{\omega}_c$ are the closures of the patches containing the interface γ . Thus $\hat{\mathbf{A}}_{22}$ is a 4×4 matrix and the associated stiffness matrix is smaller than that of the Geometric GFEM, but slightly bigger than that of the Topological GFEM.

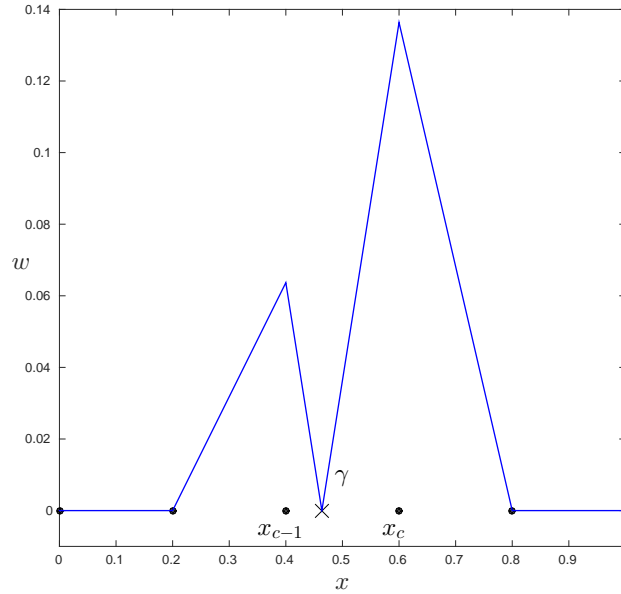


Figure 1: Enrichment function used for M-GFEM.

The Topological GFEM and M-GFEM are local in the sense that only nodes close to the interface $x = \gamma$ are enriched and that their cardinality does not change with h . We mention that another local GFEM was introduced in [20], referred to as the Corrected XFEM that employed a “Ramp cutoff function.” However the use of ramp-function yields the enrichment function w , which is quadratic in the elements $[x_{c-2}, x_{c-1}]$, $[x_c, x_{c+1}]$; the enrichment in M-GFEM is piecewise linear in every element. In contrast, the Geometric GFEM is global in the sense that $\text{card}\{x_i\}_{i \in \mathcal{R}^h} = O(h^{-1})$; for R large enough, we actually have $\mathcal{R}^h = \mathcal{N}^h$, i.e., every node of FE mesh \mathcal{T}_h could be enriched in

the Geometric GFEM for a large choice of R .

We further note that if the interface $x = \gamma$ is at one of the nodes x_i , then no enrichment is used, i.e., standard FEM can be used. However, if $x = \gamma$ is close to a node, the round-off error may create a serious problem in GFEM with the enrichments described above. In such a situation there has to be a safety check and it is advisable not to use the enrichment; we only use S_{FEM}^h . We note that not using the S_{ENR}^h part may result in loss of accuracy in the approximation, but this loss is of much smaller scale compared to the round-off error that may result if we used S_{ENR}^h . For M-GFEM, we have observed that not using any enrichment when $\min\{\gamma - x_{c-1}, x_c - \gamma\} \leq 10^{-14}h$ does not affect the accuracy of the computed solution u_h . The factor 10^{-14} in the safety-check depends on the machine precision of the computer.

First we highlight the performance of various GFEMs with respect to their accuracy. Let $\epsilon^h := \|u - u_h\|_{\mathcal{E}}$, where u_h is the computed solution using one of the GFEMs. For the Topological GFEM, one can show that $\epsilon^h \leq Ch^{1/2}$ – similar to the well-known result for FEM with uniform mesh and $\gamma \in \mathring{\tau}_c$. But for the Geometric GFEM and the M-GFEM, the rate of convergence is higher, namely, $\epsilon^h \leq Ch$. The order of convergence for Geometric GFEM and M-GFEM applied to an interface problem is thus the same as the one for the FEM applied to a problem with smooth solution.

However, the conditioning of a GFEM could be much worse than the conditioning of the FEM. Consequently, solving the linear system (3.6) could be extremely difficult. In fact, the condition number $\kappa_2(\mathbf{A})$ for the GFEM depends on the “angle” between the spaces S_{FEM}^h and S_{ENR}^h , which will be defined precisely in the next section. It has been shown in [2, 46] that if the angles between the spaces S_{FEM}^h and S_{ENR}^h are “not too small,” then $\kappa_2(\mathbf{A}) = O(h^{-2}) = \kappa_2(\mathbf{A}_{11})$, i.e., the conditioning of the GFEM is not worse than that of the standard FEM. It was shown in [2] that the “angle” could become very small for typical enrichments of GFEM used in practice, and $\kappa_2(\mathbf{A})$ could be $O(h^{-4})$.

Therefore, to design a well-conditioned GFEM, one has to choose the enrichment function w and the enrichment space S_{ENR}^h such that the “angle” between S_{FEM}^h and S_{ENR}^h “is not too small”, i.e., stays bounded away from 0.

Stable GFEM (SGFEM): A GFEM is called an SGFEM if (a) $\epsilon^h \leq Ch$ and (b) “angle” between S_{FEM}^h and S_{ENR}^h stays bounded away from 0 for all h . Specifically, for the interface problem (3.1), we let $w^* = |x - \gamma|$. S_{ENR}^h is defined with the enrichment $w = w^* - \mathcal{I}_h w^*$, where $\mathcal{I}_h w^*$ is the piecewise linear interpolant of w^* with respect to the triangulation \mathcal{T}_h . Clearly w is continuous in Ω and $w = 0$ outside τ_c . Again we consider $\mathcal{R}^h = \{x_{c-1}, x_c\}$, and consequently $\hat{\mathbf{A}}_{22}$ is a 2×2 matrix. The GFEM with enrichment w defined above yields $\epsilon^h \leq Ch$, the angle between the associated S_{ENR}^h and S_{FEM}^h is bounded away from 0, and $\kappa_2(\mathbf{A}) = O(h^{-2})$. Thus this GFEM is indeed an SGFEM. It is local, and *will be referred to as the SGFEM in this paper*. This enrichment and the associated GFEM was used in [2, 34]. Note that the same procedure applied to

the enrichment in Corrected XFEM will yield $\mathcal{R}^h = \{x_{c-2}, x_{c-1}, x_c, x_{c+1}\}$, i.e., Corrected XFEM will require more degrees of freedom.

We mention that M-GFEM is well-conditioned. But in higher dimensions, the conditioning of M-GFEM is not robust with respect to the position of the interface Γ with respect to the mesh, which we will show in the next section. On the other hand, the Geometric GFEM is badly conditioned in the sense that $\kappa_2(\mathbf{A}) = O(h^{-4})$. Note however that one has to choose h small enough, depending on R , to see this effect.

In summary, we say that the Topological GFEM is not as accurate as other GFEMs. The Geometric GFEM, though accurate, is not well conditioned for all h . The M-GFEM is accurate but the conditioning is not robust in higher dimensions. However, the SGFEM is accurate as well as robustly well-conditioned. These features will be shown for 2-D problems in the later sections of this paper.

4 Straight interface problem

In this section, we discuss the GFEM applied to a 2-D problem with a straight interface. We consider a specific manufactured problem such that the solution does not have any singularities. This problem will allow us to easily show the process of extending the 1-D ideas, presented in Section 3, to 2-D problems without the technicalities involved in a general interface problem. This problem will allow us to compare the robustness of various GFEMs considered in this paper. Moreover the straight interface problem could actually be viewed as a “laboratory problem.”

Consider the domain $\Omega = (0, 1) \times (0, 1)$. For a given $d_0 > 0$ and $0 < \theta_0 < \tan^{-1}(1/d_0)$, let $\Gamma := \{\mathbf{x} \in \Omega : \gamma(\mathbf{x}) = 0\}$ be the straight interface, where $\gamma(\mathbf{x}) = 0$ is the straight-line passing through the point $A(-d_0, 1)$ with slope $\mathcal{M} = -\tan(\theta_0)$, as shown in Figure 2, where we have chosen $d_0 = 1 - 1/\sqrt{2}$ and $\theta_0 = \pi/6$. We set $\Omega_0 := \Omega \cap \{\mathbf{x} : \gamma(\mathbf{x}) < 0\}$ and $\Omega_1 := \Omega \cap \{\mathbf{x} : \gamma(\mathbf{x}) > 0\}$. Note that varying d_0 and θ_0 , we can vary the interface Γ .

We consider the problem (2.4) with $f \equiv 0$, $q \equiv 0$, and g_N satisfying the compatibility condition (2.5). Moreover, we consider $a_i(\mathbf{x}) = a_i$, $i = 0, 1$, where a_0, a_1 are strictly positive constants. The solution $u \in \mathcal{E} = H^1(\Omega)$ of (2.4) exists and is unique up to an additive constant. We set $\|u\|_{\mathcal{E}} := B(u, u)^{1/2}$.

For the computations presented in this section, we will consider the manufactured solution of (2.4), given by

$$u_{ex} = \begin{cases} A_0 r^\alpha \cos[\alpha(\theta - \theta_0)] + B_0 r^\alpha \sin[\alpha(\theta - \theta_0)] + C, & \theta \leq \theta_0, \\ A_1 r^\alpha \cos[\alpha(\theta - \theta_0)] + B_1 r^\alpha \sin[\alpha(\theta - \theta_0)] + C, & \theta \geq \theta_0, \end{cases} \quad (4.1)$$

where (r, θ) is the polar coordinate centered at A and the polar line $\{\mathbf{x} = (x_1, x_2) : -d_0 < x_1, x_2 = 1\}$. We choose A_0, A_1, B_0, B_1 such that u_{ex} is continuous in Ω and $a(\mathbf{x}) \frac{\partial u_{ex}}{\partial n}$ is continuous across

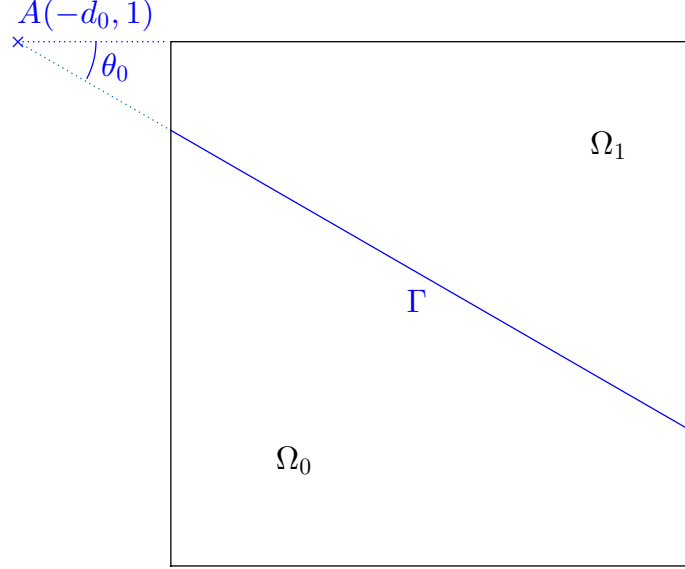


Figure 2: Straight interface problem.

the interface $\Gamma(\theta = \theta_0)$. Then, we choose C such that $u_{ex}(0, 0) = 0$. We also consider $\alpha \geq 1$. Clearly, u_{ex} is continuous in Ω with no singularity in $\overline{\Omega}$. Also $f = 0$ and g_N , in (2.4), is obtained from u_{ex} using (2.3).

We now describe the GFEM in 2-D. Let \mathcal{T}_h be a uniform finite element triangulation of Ω with nodes $\mathbf{x}_i = (i_1 h, i_2 h)$, where for a given positive integer m , we define $h = \frac{1}{m}$ and $\mathbf{i} \in \mathcal{N}^h := \{(i_1, i_2) : i_1, i_2 = 0, 1, 2, \dots, m\}$. We denote the set of elements τ , which are closed triangles with nodes as their vertices, by E . For each node \mathbf{x}_i , we define $\overline{\omega}_i = \{\cup \tau : \mathbf{x}_i \text{ is a vertex of } \tau \in E\}$. For the given triangulation \mathcal{T}_h , $\overline{\omega}_i$ is the union of 1, 2, 3 or 6 elements with the vertex \mathbf{x}_i depending on its position in Ω . The open set ω_i is the patch associated with the node \mathbf{x}_i . It is clear that $\Omega = \cup_{\mathbf{i} \in \mathcal{N}^h} \omega_i$. Let N_i be the usual piecewise linear hat-function associated with node \mathbf{x}_i with $\text{supp}\{N_i\} = \overline{\omega}_i$ and $N_i(\mathbf{x}_i) = 1$. We set

$$E_\Gamma := \{\tau \in E : \tau \cap \Gamma \neq \emptyset\}.$$

The approximation space S^h of the GFEM in 2-D is constructed similarly to the description given in Section 3. We set

$$S_{FEM}^h = \text{span}\{N_{\mathbf{i}}, \mathbf{i} \in \mathcal{N}_d^h\}, \quad \text{where } \mathcal{N}_d^h = \{\mathbf{i} \in \mathcal{N}^h : \mathbf{i} \neq (0, 0)\}.$$

The functions in the space S_{FEM}^h vanish at the node $(0, 0)$ and $\dim\{S_{FEM}^h\} = (m+1)^2 - 1$. For a given enrichment function w that mimics the exact solution, we also define the *enrichment space*

$$S_{ENR}^h = \text{span}\{wN_{\mathbf{i}}, \mathbf{i} \in \mathcal{R}^h \subset \mathcal{N}^h\}.$$

The particular choice of the enrichment function w and of the set of indices $\mathcal{R}^h \subset \mathcal{N}^h$ defines a distinct GFEM. The set $\{\mathbf{x}_{\mathbf{i}}\}_{\mathbf{i} \in \mathcal{R}^h}$ denotes the set of *enriched nodes*. The approximation space $S^h \subset \mathcal{E}$ of GFEM is given by (3.3), namely, $S^h = S_{FEM}^h \oplus S_{ENR}^h$.

The GFEM solution $u_h \in S^h$ satisfies the finite dimensional problem (3.2) with $F(v) = \int_{\partial\Omega} g \nabla v dx$ and is given in the form $u_h = \sum_{\mathbf{i} \in \mathcal{N}_d^h} c_{1,\mathbf{i}} N_{\mathbf{i}} + \sum_{\mathbf{k} \in \mathcal{R}^h} c_{2,\mathbf{k}} w N_{\mathbf{k}}$, where $c_{1,\mathbf{i}}, c_{2,\mathbf{k}}$ is the solution of the linear system (3.5). Note that we solve the diagonally scaled linear system (3.6), instead of (3.5). Note also that if $w \equiv 0$, i.e., if no enrichment is used, $S^h = S_{FEM}^h$ and the GFEM is the standard FEM.

We will now describe various GFEMs based on the particular choices of the enrichment function w and the set of indices \mathcal{R}^h for the enriched nodes. The enrichment function w for the interface problem is based on the so-called *distance function*

$$w^*(\mathbf{x}) := \text{dist}(\mathbf{x}, \Gamma).$$

Note that $w^*(\mathbf{x})$ is continuous in Ω , it is linear in Ω_0 and Ω_1 , and $w^*(\mathbf{x}) = 0$ for $\mathbf{x} \in \Gamma$.

Geometric GFEM: The enrichment space S_{ENR}^h is constructed with the enrichment function $w(\mathbf{x}) = w^*(\mathbf{x})$, whereas the set of indices \mathcal{R}^h is given by $\mathcal{R}^h = \{\mathbf{i} \in \mathcal{N}^h : \text{dist}(\mathbf{x}_{\mathbf{i}}, \Gamma) \leq R\}$ for a fixed R , independent of h . Note that unlike in Section 3, $\text{card}\{\mathbf{x}_{\mathbf{i}}\}_{\mathbf{i} \in \mathcal{R}^h} = O(h^{-2})$. The set of enriched nodes thus contains all the nodes within a fixed distance R from the interface.

Topological GFEM: The same enrichment function $w(\mathbf{x}) = w^*(\mathbf{x})$ is used to define S_{ENR}^h as in the Geometric GFEM. Here however, we use $\mathcal{R}^h = \{\mathbf{i} \in \mathcal{N}^h : \omega_{\mathbf{i}} \cap \Gamma \neq \emptyset\}$. Note that the set of enriched nodes $\{\mathbf{x}_{\mathbf{i}} : \mathbf{i} \in \mathcal{R}^h\}$ is the union of all the vertices of the elements $\tau \in E_{\Gamma}$. Again unlike in Section 3, $\text{card}\{\mathbf{x}_{\mathbf{i}}\}_{\mathbf{i} \in \mathcal{R}^h} = O(h^{-1})$.

M-GFEM: The enrichment function used in M-GFEM is slightly different than w used in Geometrical or Topological GFEM and is given by

$$w(\mathbf{x}) = \begin{cases} w^*(\mathbf{x}), & \mathbf{x} \in \tau \in E_{\Gamma}, \\ \text{linear function}, & \mathbf{x} \in \tau \in E \setminus E_{\Gamma}, \end{cases}$$

$$w(\mathbf{x}_{\mathbf{i}}) = 0, \quad \mathbf{x}_{\mathbf{i}} \text{ is not a vertex of } \tau \in E_{\Gamma}.$$

Furthermore, the indices of enriched nodes are given by

$$\mathcal{R}^h = \{\mathbf{i} \in \mathcal{N}^h : \mathbf{x}_{\mathbf{i}} \in \bigcup_{\omega_j \cap \Gamma \neq \emptyset} \bar{\omega}_j\}.$$

Note that unlike in Section 3, $\text{card}\{\mathbf{x}_{\mathbf{i}}\}_{\mathbf{i} \in \mathcal{R}^h} = O(h^{-1})$.

The enriched nodes $\mathbf{x}_{\mathbf{i}}, \mathbf{i} \in \mathcal{R}^h$ are shown in Figure 3, where $h = 1/16$, $d_0 = 1 - 1/\sqrt{2}$, $\theta_0 = \pi/6$ and $R = 1/3$.

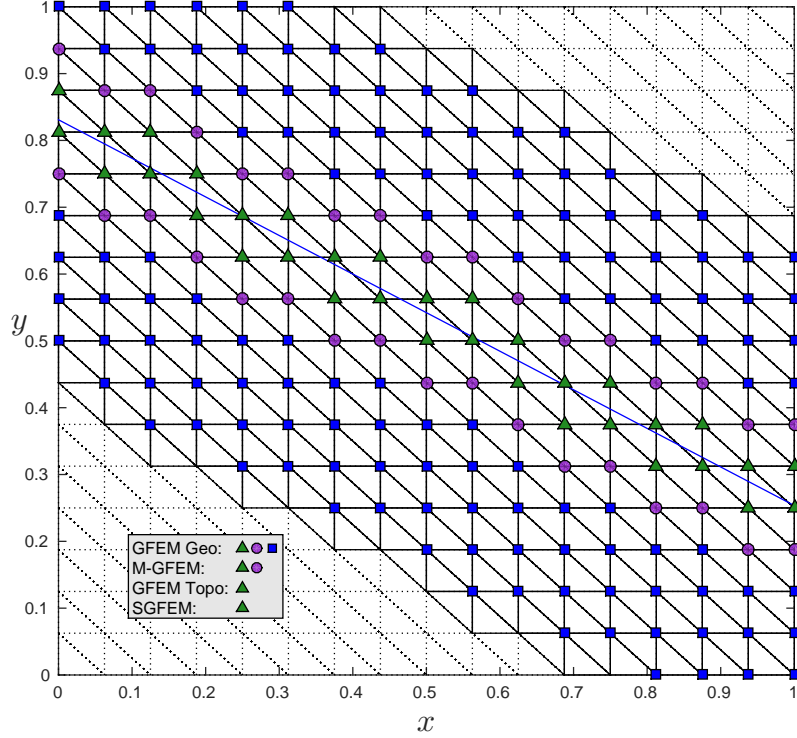


Figure 3: Nodes enriched for the straight interface problem.

First, we consider the exact solution $u \in \mathcal{E}$ of (2.4) given in (4.1) with $d_0 = 1 - 1/\sqrt{2}$, $\theta_0 = \pi/6$, $a_0 = 1$ and $a_1 = 10$. Note that the interface Γ is not aligned with the mesh. We computed the error $\|u - u_h\|_{\mathcal{E}}$, where u_h is the solution of (3.2) associated with the GFEMs described above (we chose $R = 1/3$ for Geometric GFEM). The log-log plot of the (relative) error is given in Figure 4. It is clear that Geometric GFEM and M-GFEM yield the convergence of $O(h)$, whereas the order of convergence for the Topological GFEM is only $O(h^{1/2})$. This suboptimal order of convergence for Topological GFEM has been reported in the literature [28, 40, 44].

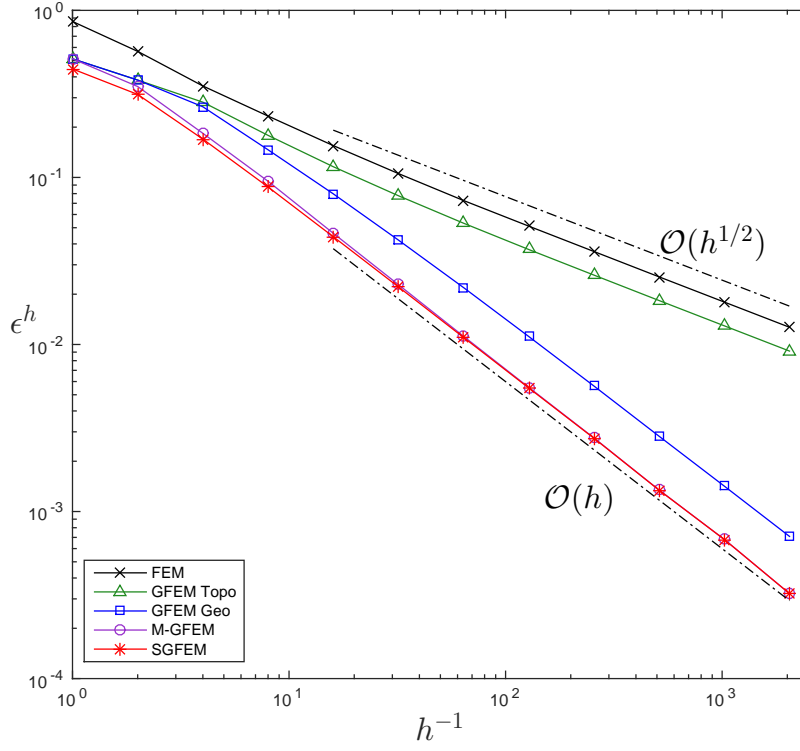


Figure 4: Relative error in the energy norm against h .

The approximation property of GFEM, as described in the last paragraph, is undoubtedly very important. However, it is equally important that the GFEM be well-conditioned in order for the linear system (3.6) to be solved efficiently. Towards this end, we first compute the condition number of the scaled stiffness matrix \mathbf{A} in (3.6) for different values of h . The results are collected in Figure 5 (solid lines), again with $d_0 = 1 - 1/\sqrt{2}$, $\theta_0 = \pi/6$, $R = 1/3$, $a_0 = 1$ and $a_1 = 10$. It is clear that for the Topological GFEM and M-GFEM, the condition number $\kappa_2(\mathbf{A}) = O(h^{-2})$. We mention that $\kappa_2(\mathbf{A}_{11}) = O(h^{-2})$, where \mathbf{A}_{11} is the stiffness matrix of the standard FEM. In other words, the conditioning of the Topological GFEM and the M-GFEM is of the same order as that of a standard FEM. On the other hand, $\kappa_2(\mathbf{A}) = O(h^{-4})$ for the Geometrical GFEM, which is much worse than the Topological GFEM and the M-GFEM. We will show later in this paper that conditioning also plays a major role in the iterative solution of (3.6). We also mention that $\kappa_2(\mathbf{A}_{22})$ is bounded for all considered GFEMs, except for Geometric GFEM for which we have $\kappa_2(\mathbf{A}_{22}) = O(h^{-2})$ (see dashed lines in Figure 5).

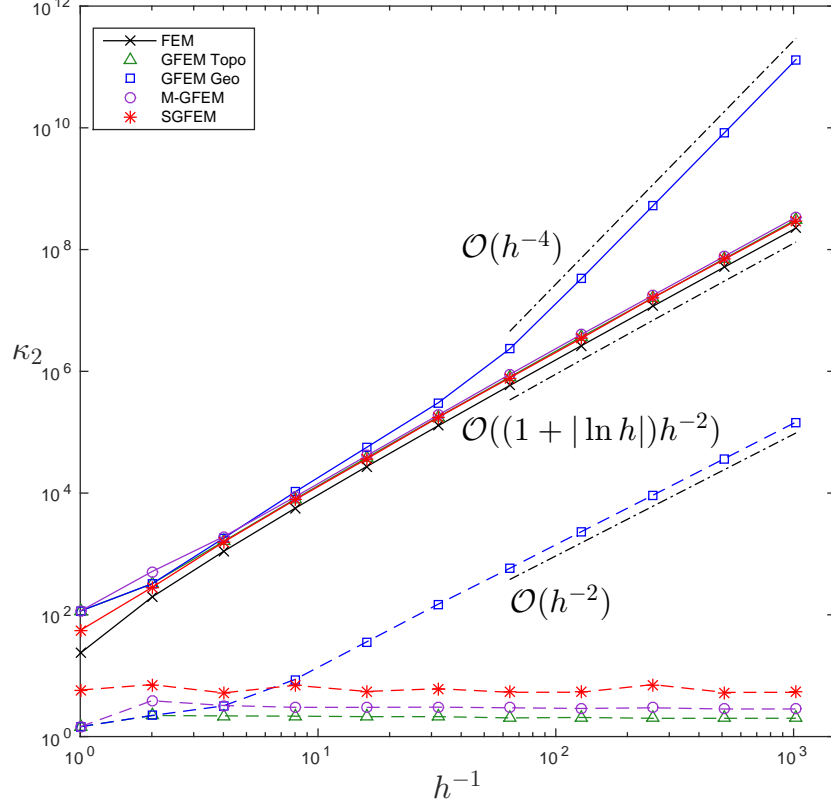


Figure 5: Condition number of the scaled stiffness matrices \mathbf{A} (solid lines) and \mathbf{A}_{22} (dashed lines) against h .

The conditioning of GFEM depends on the “angle” between the spaces S_{FEM}^h and S_{ENR}^h , which is characterized by the quantity

$$\cos(\vartheta(S_{FEM}^h, S_{ENR}^h)) := \max_{u \in S_{FEM}^h, v \in S_{ENR}^h} \frac{|B(u, v)|}{\|u\|_{\mathcal{E}} \|v\|_{\mathcal{E}}},$$

where $\vartheta(S_{FEM}^h, S_{ENR}^h) \in [0, 90^\circ]$ could be interpreted as the smallest angle between S_{FEM}^h and S_{ENR}^h and depends on the specific enrichment function $w(\mathbf{x})$ used in the GFEM. In particular, if $\vartheta(S_{FEM}^h, S_{ENR}^h) = 90^\circ$, then we could write $S^h = S_{FEM}^h \oplus^\perp S_{ENR}^h$, i.e., the two spaces are in orthogonal direct sum. On the other hand, if $\vartheta(S_{FEM}^h, S_{ENR}^h) = 0$, then we simply have $S^h = S_{FEM}^h + S_{ENR}^h$, and the sum is no longer direct. When inbetween, $\vartheta(S_{FEM}^h, S_{ENR}^h) \in (0, 90^\circ)$, we have $S^h = S_{FEM}^h \oplus S_{ENR}^h$, and the sum is direct but not orthogonal. It has been proved in

[2, 46] that if the angles between the spaces S_{FEM}^h and S_{ENR}^h are “not too small”, i.e., if there exist positive constants C, C_1, C_2 such that

$$\vartheta(S_{FEM}^h, S_{ENR}^h) \geq C > 0, \quad (4.2)$$

and if

$$C_1 \leq \kappa_2(\mathbf{A}_{22}) \leq C_2, \quad (4.3)$$

where \mathbf{A}_{22} as in (3.6), then

$$\kappa_2(\mathbf{A}) = O(h^{-2}) = \kappa_2(\mathbf{A}_{11}), \quad (4.4)$$

i.e., the conditioning of GFEM, with S_{FEM}^h and S_{ENR}^h satisfying the above conditions, is not worse than that of the standard FEM. Note however that conditions (4.2)–(4.3) are sufficient conditions for (4.4), i.e., they guarantee the well-conditioning of the GFEM. We further note that (4.4) holds even when the condition (4.3) is replaced by $\kappa_2(\mathbf{A}_{22}) = O(h^{-2})$. Moreover, since the functions in S_{FEM}^h vanish only at the node $(0,0)$, it is theoretically known [11] that $\kappa_2(\mathbf{A}_{11}) = O[h^{-2}(1 + |\ln h|)]$. In this paper, we will not consider the factor $|\ln h|$, and consider instead $\kappa_2(\mathbf{A}_{11}) = O(h^{-2})$, as done in (4.4). The results in [2, 46] are quite general. As long as S_{ENR}^h associated with *any chosen enrichment function satisfies the conditions (4.2)–(4.3), the well-conditioning of the associated GFEM is guaranteed*. This is not only true for interface problems, but for any problem where GFEM is used. In fact, the conditions (4.2)–(4.3) may help to construct enrichments leading to well-conditioned GFEM.

To investigate the dependence of $\kappa_2(\mathbf{A})$ on the angle between S_{FEM}^h and S_{ENR}^h , we computed the angle for Geometric GFEM, Topological GFEM, and the M-GFEM for different values of h . The results are displayed in Figure 6, again with $d_0 = 1 - 1/\sqrt{2}$, $\theta_0 = \pi/6$, $R = 1/3$, $a_0 = 1$ and $a_1 = 10$. The angle could be obtained by solving a generalized eigenvalue problem that we present in Appendix A. It is clear from Figure 6 that the angle for the Topological GFEM and M-GFEM remain bounded away from 0 for all the values of h , thus illuminating the result presented above since we have seen in Figure 5 that $\kappa_2(\mathbf{A}) = O(h^{-2})$ for the Topological GFEM and M-GFEM. Figure 6 also shows that the angle for the Geometric GFEM approaches 0 as h gets smaller. Moreover, Figure 5 indicates that $\kappa_2(\mathbf{A}) = O(h^{-4}) \gg O(h^{-2})$. This suggests that the condition (4.2) could be a necessary condition for (4.4).

Therefore, well-conditioning (of the same order as the standard FEM) of a system could be guaranteed if the GFEM uses enrichments that yield accurate approximation and if the angle between S_{FEM}^h and S_{ENR}^h is uniformly bounded away from 0.

Stable GFEM (SGFEM): A GFEM is called an SGFEM if (a) it yields the optimal rate of convergence, and (b) it satisfies the conditions (4.2)–(4.3). For the straight interface problems, a special enrichment could be obtained by a simple modification of the distance function $w^*(\mathbf{x})$ and

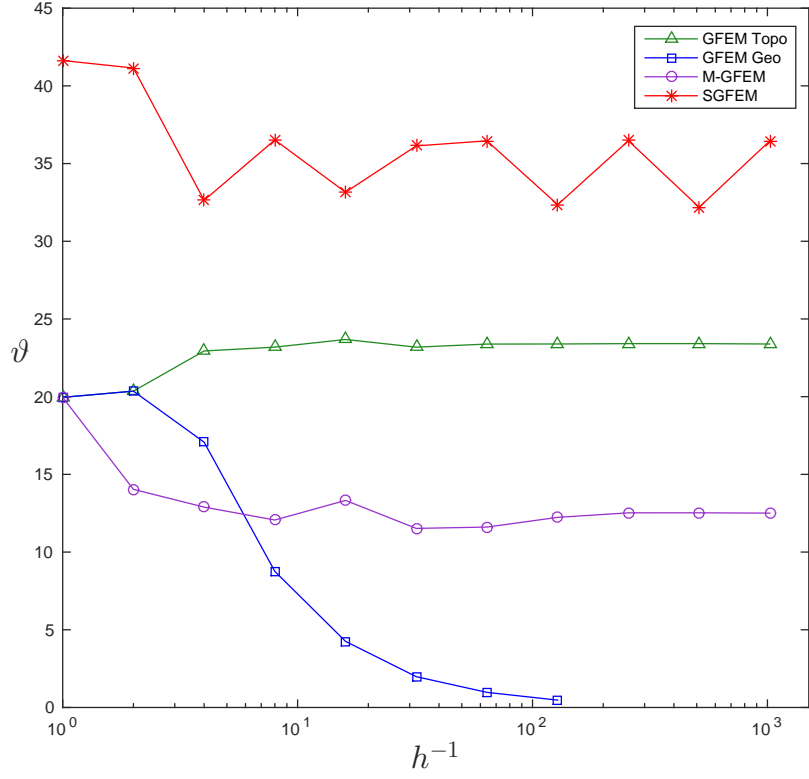


Figure 6: Angle (in degrees) between the spaces S_{FEM}^h and S_{ENR}^h against h .

the associated GFEM is indeed an SGFEM. The enrichment is defined as

$$w(\mathbf{x}) = w^*(\mathbf{x}) - \mathcal{I}_h w^*(\mathbf{x}), \quad (4.5)$$

where $\mathcal{I}_h w^*(\mathbf{x})$ is the piecewise linear interpolant of $w^*(\mathbf{x})$. Since $w^*(\mathbf{x})$ is linear on Ω_0 and Ω_1 , it is easy to see that

$$\text{supp}\{w(\mathbf{x})\} = \bigcup_{\tau \in E_\Gamma} \tau.$$

Moreover, the indices of the enriched nodes are

$$\mathcal{R}^h = \{\mathbf{i} \in \mathcal{N}^h : \omega_{\mathbf{i}} \cap \Gamma \neq \emptyset\}, \quad (4.6)$$

which is the same as the \mathcal{R}^h used in Topological GFEM (the enriched nodes are shown in Figure 3).

In Figures 4, 5 and 6, we have plotted the error $\|u - u_h\|_{\mathcal{E}}$, the condition number $\kappa_2(\mathbf{A})$ and the angle between S_{FEM}^h and S_{ENR}^h , with respect to h , for the GFEM with enrichment given by (4.5) and \mathcal{R}^h as in (4.6). It is clear that the method yields the optimal order of convergence, i.e., $\|u - u_h\|_{\mathcal{E}} = O(h)$, $\kappa_2(\mathbf{A}) = O(h^{-2})$, and the angle between S_{FEM}^h and S_{ENR}^h is bounded away from 0 for all the values of h considered in the experiment. The condition (4.2) is thus satisfied. We have also checked (see Figure 5, dashed lines) that (4.3) is satisfied. Thus the GFEM with enrichment given in (4.5)–(4.6) is indeed an SGFEM; *we will refer to this GFEM as SGFEM in the rest of this paper*. We further note in Figure 6 that the angle between S_{FEM}^h and S_{ENR}^h for the SGFEM is larger than that of the M-GFEM – this feature is central to the iterative solution of (3.6), which we will show later in this paper.

Remark 4.1 Note that the enrichment given in (4.5) and the set of enriched nodes indexed by \mathcal{R}^h in (4.6) was introduced in [34]. However, the conditioning of the GFEM or the angle between S_{FEM}^h and S_{ENR}^h was not discussed there.

Remark 4.2 It is important to note that modifying an enrichment by subtracting the piecewise linear interpolant, as we have done in (4.5), may not yield an SGFEM for other problems. It has been shown in [23] that for the crack propagation problems, modification of enrichment by subtracting a linear interpolant may yield inaccurate solutions. An additional modification of the enrichment function is needed to yield accurate solutions for such problems. However, the modification is certainly successful for the interface problems, as those considered in this paper.

It has been reported in the literature [21] that GFEM may become ill-conditioned if the interface Γ is close to the mesh lines. To investigate this problem, we consider the manufactured solution (4.1) with $\theta_0 = \pi/4$. In this case, the interface Γ is parallel to some of the mesh lines associated with the triangulation \mathcal{T}_h . We control the distance of Γ to the mesh line by controlling the parameter d_0 (see Figure 2). We have fixed $h = 1/16$ and have plotted the condition number $\kappa_2(\mathbf{A})$ for Topological GFEM, Geometric GFEM, M-GFEM, and SGFEM, as the interface Γ gets closer to the mesh line in Figure 7. In this scenario, the other parameters are $R = 1/6$, $a_0 = 1$ and $a_1 = 10$. It is clear that $\kappa_2(\mathbf{A})$ for M-GFEM “blows-up” as Γ gets closer to the mesh line; there is no appreciable change in $\kappa_2(\mathbf{A})$ for other GFEMs considered here. We mention that $\kappa_2(\mathbf{A}_{22})$ stays bounded for all GFEMs considered, even when the interface is relatively close to the mesh lines, and even for M-GFEM (not shown in Figure 7). In Figure 8, we have plotted the angle between S_{FEM}^h and S_{ENR}^h for the fixed $h = 1/16$. We clearly see that the angle for M-GFEM goes to 0 as Γ gets closer to the mesh line; angles for other GFEMs approach different but fixed values, bounded away from 0. This shows that the *conditioning of M-GFEM is not robust with respect to the position of the interface to the edges of the mesh*. We mention that in a forthcoming paper, we will prove that the GFEM with enrichment (4.5) satisfies the conditions (4.2)–(4.3), where the constants are

independent of h and of the position of Γ . In other words, *SGFEM is robust with respect to the position of the interface to the edges of the mesh*. This is also illuminated in Figure 8 where we see that the angle between S_{FEM}^h and S_{ENR}^h for the SGFEM approaches a fixed value bounded away from 0. However, similarly to the 1-D interface problem, there has to be a safety-check and it is advisable not to enrich a node if the interface is very close to it; otherwise, round-off errors could contaminate the solution.

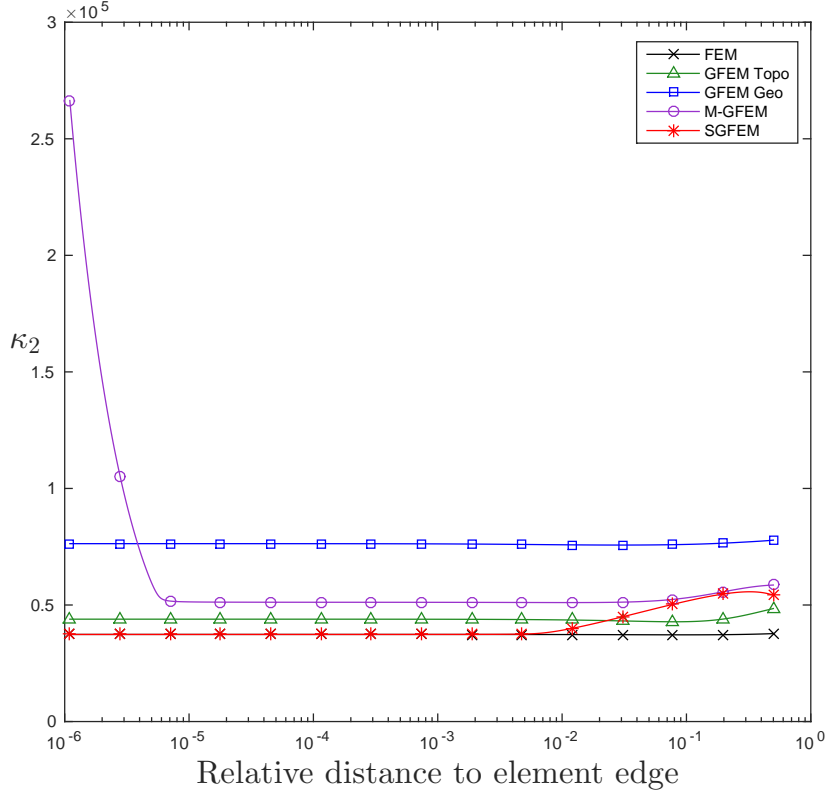


Figure 7: Condition number of the scaled stiffness matrix \mathbf{A} when the interface gets closer to the mesh lines.

We summarize the results mentioned above about different GFEMs in the following table.

	FEM	GFEM (Topo)	GFEM (Geo)	GFEM (M)	SGFEM
Order of Conv.	$O(h^{1/2})$	$O(h^{1/2})$	$O(h)$	$O(h)$	$O(h)$
Angle		bounded away from 0	$\rightarrow 0$ as $h \rightarrow 0$	$\rightarrow 0$ only as $\Gamma \rightarrow \text{edge}$	bounded away from 0
$\kappa_2(\mathbf{A})$	$O(h^{-2})$	$O(h^{-2})$	$O(h^{-4})$	$O(h^{-2})$	$O(h^{-2})$
Robustness	yes	yes	yes	no	yes

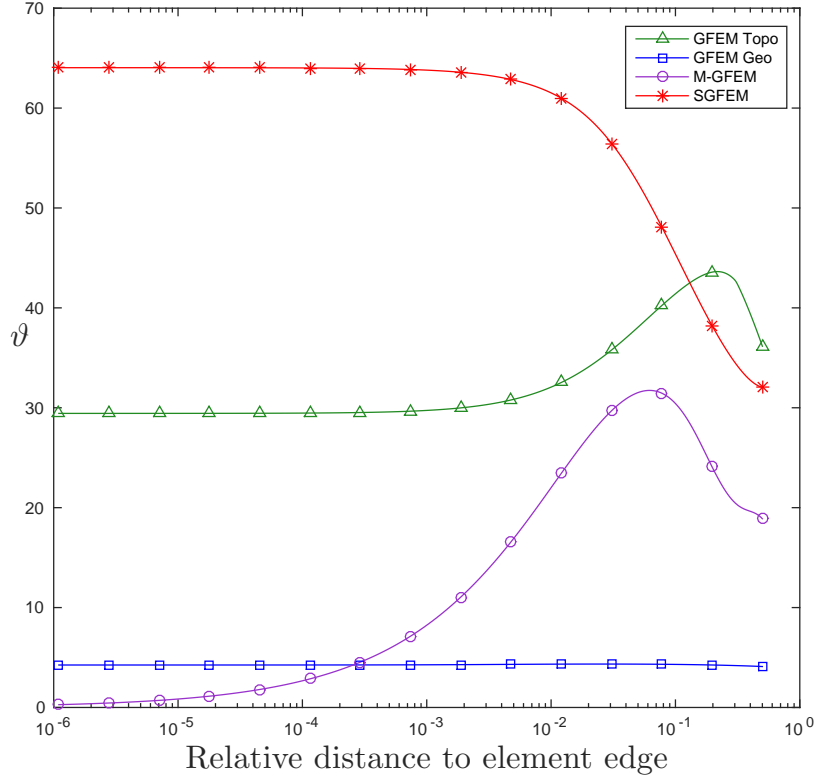


Figure 8: Angle (in degrees) between the spaces S_{FEM}^h and S_{ENR}^h when the interface gets closer to the mesh lines.

We thus conclude that among all the GFEMs considered in this section, the SGFEM is the only method that has all the desired features – it yields accurate approximation, it is well-conditioned, and it is robust.

5 Circular interface problem

In this section, we discuss the GFEM applied to a 2-D problem with a circular interface. We consider a specific manufactured problem such that the solution does not have any singularities. This problem will be solved by extending what was done in Section 4 on a straight interface to a circular interface.

Consider the domain $\Omega = (0, 1) \times (0, 1)$. For given $r_c > 0$ and $(x_c, y_c) \in \Omega$, let $\Gamma := \{\mathbf{x} \in \Omega : \gamma(\mathbf{x}) = 0\}$ be the interface, where $\gamma(\mathbf{x}) = (x - x_c)^2 + (y - y_c)^2 - r_c^2$ is the circle of center $A(x_c, y_c)$

and radius r_c , as shown in Figure 9, where we have chosen $r_c = 1/\sqrt{10}$ and $(x_c, y_c) = (1/\sqrt{5}, 1/\sqrt{3})$. We set $\Omega_0 := \Omega \cap \{\mathbf{x} : \gamma(\mathbf{x}) < 0\}$ and $\Omega_1 := \Omega \cap \{\mathbf{x} : \gamma(\mathbf{x}) > 0\}$. Note that when varying x_c, y_c, r_c , the interface Γ varies as well.

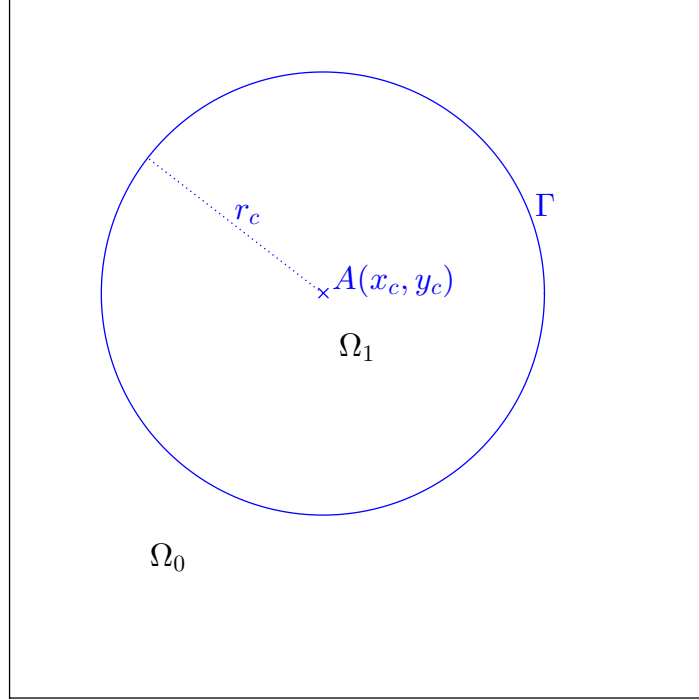


Figure 9: Circular interface problem.

We consider the problem (2.4) with $f \equiv 0$, $q \equiv 0$, and g_N satisfying the compatibility condition (2.5). Moreover, we consider $a_i(\mathbf{x}) = a_i$, $i = 0, 1$, where a_0, a_1 are strictly positive constants. The solution $u \in \mathcal{E} = H^1(\Omega)$ of (2.4) exists and is unique up to an additive constant. We set $\|u\|_{\mathcal{E}} := B(u, u)^{1/2}$.

For the computations presented in this section, we will consider the manufactured solution of (2.4), given by

$$u_{ex} = \begin{cases} r^2 \cos(2\theta) + C, & r \leq r_c, \\ B_0 r^2 \cos(2\theta) + B_1 r^{-2} \cos(2\theta) + C, & r \geq r_c, \end{cases} \quad (5.1)$$

where (r, θ) is the polar coordinate centered at A . We choose B_0, B_1 such that u_{ex} and $a(\mathbf{x}) \frac{\partial u_{ex}}{\partial n}$ are continuous across the interface Γ ($r = r_c$) and then C so that $u_{ex}(0, 0) = 0$. It is clear that u_{ex}

is continuous in Ω with no singularity in $\overline{\Omega}$. Moreover $f = 0$ and g_N , in (2.4), is obtained from u_{ex} using (2.3).

Let us first comment on FEM. We use the same discretization and notations as in Section 4 to define $\mathbf{x}_i, \tau, \omega_i, N_i$ with an understanding that they depend on h . During the assembling of the stiffness matrix \mathbf{A}_{11} we need to evaluate quantities of the form

$$B(N_j, N_i) = \int_{\Omega} a \nabla N_j \nabla N_i d\mathbf{x}.$$

With a being discontinuous on Γ , we thus have to split the integration domain into two complementary parts, one in Ω_0 and one in Ω_1 . Each of these sub-parts is then no longer polygonal as the interface is curved, numerical integration up to machine precision would thus be costly and difficult to implement. We propose another approach instead. Considering that the interface can be described by $\gamma(\mathbf{x}) = (x - x_c)^2 + (y - y_c)^2 - r_c^2$, we can readily invert this expression to find the intersections of the finite element triangulation with the interface. This gives a set of points whose convex hull forms a polygon, noted $\tilde{\Gamma}$ – see Figure 10, where $h = 1/2$, $r_c = 1/\sqrt{10}$ and $(x_c, y_c) = (1/\sqrt{5}, 1/\sqrt{3})$. We will use this polygon instead of the circular interface.

We define a “perturbed” $a(\mathbf{x})$, noted $\tilde{a}(\mathbf{x})$, whose value is a_0 outside the polygon, domain denoted $\tilde{\Omega}_0$ and a_1 inside, domain denoted $\tilde{\Omega}_1$. We also define ω as the difference $\omega := \tilde{\Omega}_0 \setminus \Omega_0$ (we could equivalently define ω as the difference $\Omega_1 \setminus \tilde{\Omega}_1$). This approach could be considered as a perturbation of the original problem. We can now see why we first studied the straight interface problem. With the perturbed interface, instead of (3.1) we now consider the perturbed variational problem: find $\tilde{u} \in \mathcal{E}$ satisfying

$$\tilde{B}(\tilde{u}, v) = F(v), \quad \text{for all } v \in \mathcal{E}, \quad (5.2)$$

where

$$\tilde{B}(u, v) := \int_{\Omega} \tilde{a} \nabla u \cdot \nabla v d\mathbf{x}.$$

Note that while the solution $u \in \mathcal{E}$ (5.1) of the original problem (3.1) does not have any singularity in $\overline{\Omega}$, the perturbed solution $\tilde{u} \in \mathcal{E}$ of the perturbed problem (5.2) may exhibit singularities due to the corners of the perturbed interface $\tilde{\Gamma}$. However, we do not solve (5.2). Instead, we solve the finite dimensional problem: find $u_h \in S^h$ satisfying

$$\tilde{B}(u_h, v) = F(v), \quad \text{for all } v \in S^h, \quad (5.3)$$

so the solution $u_h \in S^h$ does not have any singularity in $\overline{\Omega}$.

Since numerical solutions $u_h \in S^h$ of (5.3) satisfy Galerkin orthogonality only with respect to \tilde{u} and $\tilde{B}(\cdot, \cdot)$ and not with respect to u and $B(\cdot, \cdot)$, we compute the discretization error as $\epsilon^h := \left| \|u\|_{\mathcal{E}}^2 - \|u_h\|_{\tilde{\mathcal{E}}}^2 \right|^{1/2}$, where $\|v\|_{\tilde{\mathcal{E}}} := \tilde{B}(v, v)^{1/2}$. First, the energy of the exact solution $u \in \mathcal{E}$

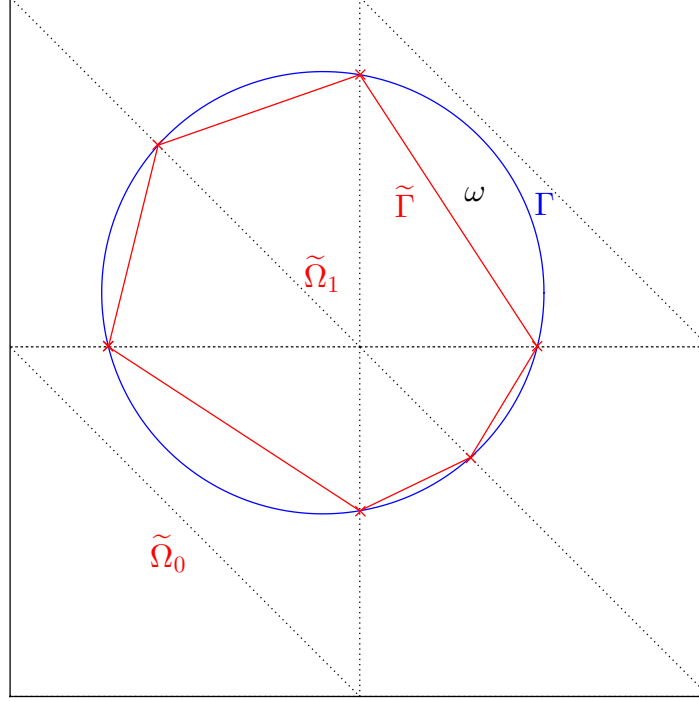


Figure 10: Circular interface and perturbed interface.

is still computed with respect to $B(\cdot, \cdot)$ (i.e., the true interface Γ), while the energy of the discrete solution $u_h \in S^h$ is now computed with respect to $\tilde{B}(\cdot, \cdot)$ (i.e., the perturbed interface $\tilde{\Gamma}$). Secondly, we compute the “difference of the energy norms” and not “energy norm of the difference” (note that the two usually coincide thanks to Galerkin orthogonality). We use this unusual definition of the discretization error for computational reasons. This definition avoids integrating the quantity $a \nabla(u - u_h) \cdot \nabla(u - u_h)$ on each element $\tau \in E$, which would be costly and difficult to implement due to the curved interface Γ . Instead, we thus compute the discretization error as $\epsilon^h = \left| \|u\|_{\mathcal{E}}^2 - \|u_h\|_{\tilde{\mathcal{E}}}^2 \right|^{1/2}$, and it holds

$$\left| \|u - u_h\|_{\mathcal{E}}^2 - \left| \|u\|_{\mathcal{E}}^2 - \|u_h\|_{\tilde{\mathcal{E}}}^2 \right| \right| \leq O(h^k),$$

where $k = 3/2$ for FEM and $k = 2$ for GFEM & SGFEM. The proof of this result can be found in Appendix B. Note that we have committed two crimes: first, the perturbation of the interface

from Γ to $\tilde{\Gamma}$, secondly, the computation of the error as $\left| \|u\|_{\mathcal{E}}^2 - \|u_h\|_{\mathcal{E}}^2 \right|^{1/2}$ and not as $\|u - u_h\|_{\mathcal{E}}$. However, these two crimes have limited effects compared to the true discretization error $\|u - u_h\|_{\mathcal{E}}$ which is of order $O(h^{1/2})$ for FEM and $O(h)$ for GFEM & SGFEM.

Finally, note that the quantity $\|u\|_{\mathcal{E}}^2$ in ϵ^h can be calculated as $\int_{\partial\Omega} u g_N ds$, thus avoiding the curved interface Γ .

Let us now describe the GFEM in 2-D. The approximation space S^h of the GFEM is once again given by $S^h = S_{FEM}^h \oplus S_{ENR}^h$.

In this section, we will only consider the so-called M-GFEM of the previous section. Recall that Topological GFEM does not recover the optimal rate of convergence in terms of error, while Geometric GFEM is badly conditioned; we thus do not discuss them anymore. For brevity, we will refer to M-GFEM simply as GFEM in this section and the next. The enrichment function w for the circular problem is again based on the so-called *distance function*. However, note that $\text{dist}(\mathbf{x}, \Gamma)$ is quadratic in Ω_0 and Ω_1 . Moreover, $\text{dist}(\mathbf{x}, \tilde{\Gamma})$ is quadratic in some regions of Ω . In order to facilitate numerical integration, we would like to use a piecewise linear enrichment function. Here is how we were able to obtain such a function. We start with $w^\diamond(\mathbf{x}) = \text{dist}(\mathbf{x}, \Gamma)$ the distance to the interface Γ . Next, we perform a triangulation of \mathcal{T}_h using $\tilde{\Gamma}$ as an edge constraint: each element in E_Γ is divided into elementary triangles whose edges do not cross $\tilde{\Gamma}$. We then compute the linear interpolant of $w^\diamond(\mathbf{x})$ on this triangulation, thus giving $w^*(\mathbf{x})$. It is continuous and piecewise linear in Ω . We mention the presence of “shadow interfaces” due to the triangulation, but the additional weak discontinuities (kinks) are controlled by the true distance to the interface rather than by arbitrary factors. Then, the enrichment function used in GFEM is given by

$$w(\mathbf{x}) = \begin{cases} w^*(\mathbf{x}), & \mathbf{x} \in \tau \in E_\Gamma, \\ \text{linear function}, & \mathbf{x} \in \tau \in E \setminus E_\Gamma, \end{cases}$$

$$w(\mathbf{x}_i) = 0, \quad \mathbf{x}_i \text{ is not a vertex of } \tau \in E_\Gamma.$$

Furthermore, the set of indices of enriched nodes is given by

$$\mathcal{R}^h = \{\mathbf{i} \in \mathcal{N}^h : \mathbf{x}_i \in \bigcup_{\omega_j \cap \Gamma \neq \emptyset} \overline{\omega}_j\}.$$

The enriched nodes $\mathbf{x}_i, \mathbf{i} \in \mathcal{R}^h$ are shown in Figure 11, where $h = 1/16$, $r_c = 1/\sqrt{10}$ and $(x_c, y_c) = (1/\sqrt{5}, 1/\sqrt{3})$. Note that like in Section 4, $\text{card}\{\mathbf{x}_i\}_{\mathbf{i} \in \mathcal{R}^h} = O(h^{-1})$.

First, let us consider the exact solution $u \in \mathcal{E}$ of (2.4) given in (5.1) with $r_c = 1/\sqrt{10}$, $(x_c, y_c) = (1/\sqrt{5}, 1/\sqrt{3})$, $a_0 = 1$ and $a_1 = 10$. We computed the error $\epsilon^h = \left| \|u\|_{\mathcal{E}}^2 - \|u_h\|_{\mathcal{E}}^2 \right|^{1/2}$, where u_h is the solution of (5.3) associated with the GFEM described above. The log-log plot of the (relative) error is given in Figure 12. It is clear that GFEM yields the convergence of $O(h)$, whereas the order of convergence for standard FEM is only $O(h^{1/2})$. This is very similar to the results of Section 4.

Concerning the condition number of the scaled stiffness matrix \mathbf{A} in (3.6) for different values of h , we display the results in Figure 13 (solid lines), again with $r_c = 1/\sqrt{10}$, $(x_c, y_c) = (1/\sqrt{5}, 1/\sqrt{3})$,

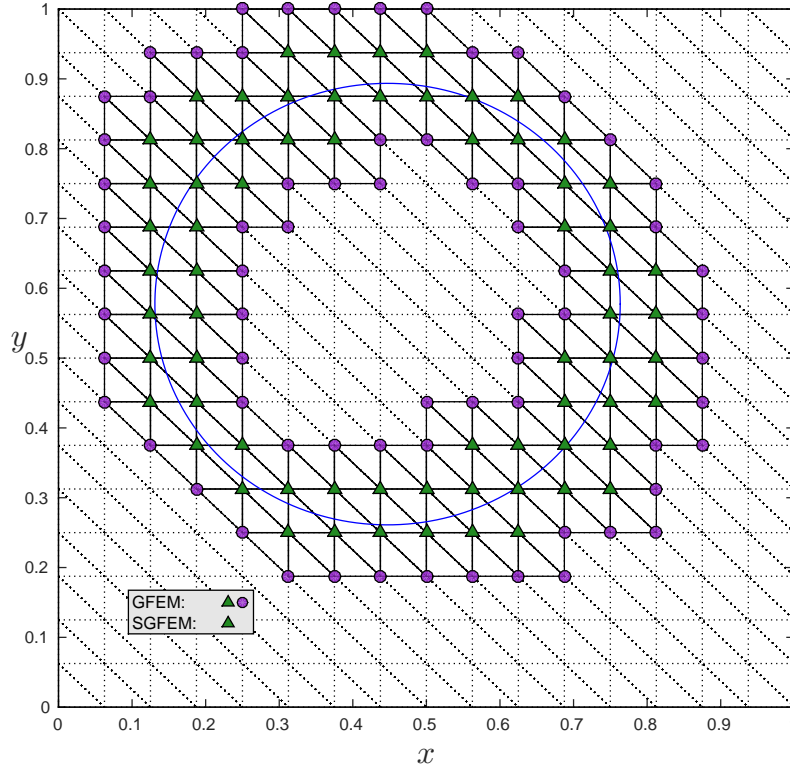


Figure 11: Nodes enriched for the circular interface problem.

$a_0 = 1$ and $a_1 = 10$. It is clear that for GFEM, the condition number $\kappa_2(\mathbf{A}) = O(h^{-2})$. We mention that $\kappa_2(\mathbf{A}_{11}) = O(h^{-2})$, where \mathbf{A}_{11} is the stiffness matrix of the standard FEM. In other words, the conditioning of GFEM is of the same order as that of a standard FEM. Again, this is very similar to the results in Section 4. We also mention that $\kappa_2(\mathbf{A}_{22})$ is bounded for all values of h (see dashed lines in Figure 13).

We have computed the angle for the GFEM for different values of h and displayed the results in Figure 14, again with $r_c = 1/\sqrt{10}$, $(x_c, y_c) = (1/\sqrt{5}, 1/\sqrt{3})$, $a_0 = 1$ and $a_1 = 10$. It is clear that the angle remains bounded away from 0 for all the values of h and, it thus illuminates the result presented above since we have seen in Figure 13 that $\kappa_2(\mathbf{A}) = O(h^{-2})$.

Stable GFEM (SGFEM): The enrichment is defined as

$$w(\mathbf{x}) = w^*(\mathbf{x}) - \mathcal{I}_h w^*(\mathbf{x}), \quad (5.4)$$

where $\mathcal{I}_h w^*(\mathbf{x})$ is the piecewise linear interpolant of $w^*(\mathbf{x})$ with respect to the triangulation \mathcal{T}_h .

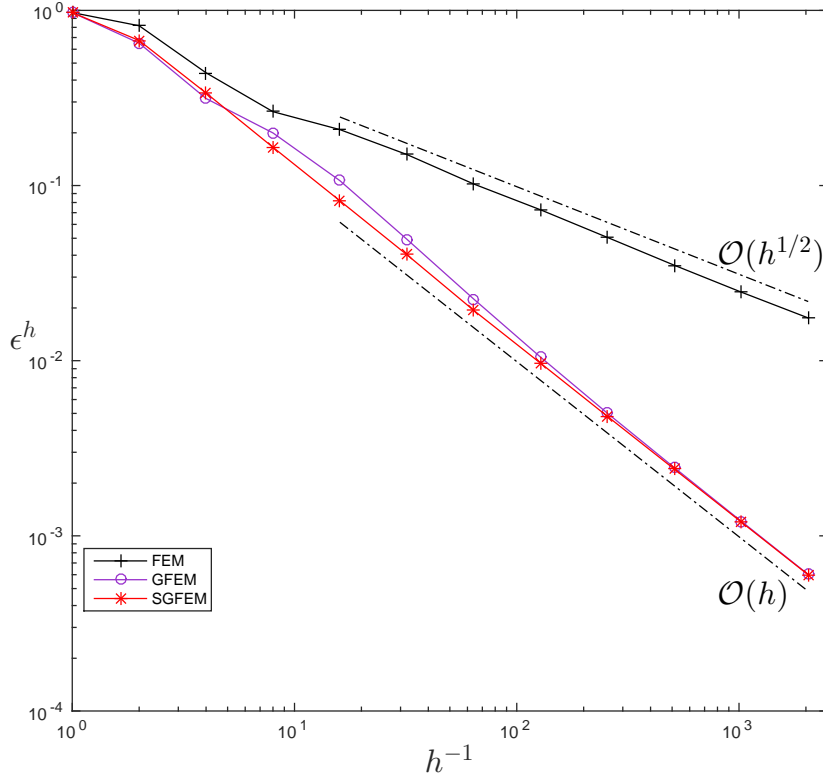


Figure 12: Relative error in the energy norm against h .

Since $w^*(\mathbf{x})$ is linear on $\tau \in E \setminus E_\Gamma$, it is easy to see that

$$\text{supp}\{w(\mathbf{x})\} = \bigcup_{\tau \in E_\Gamma} \tau.$$

Moreover, the indices of the enriched nodes are

$$\mathcal{R}^h = \{\mathbf{i} : \omega_{\mathbf{i}} \cap \Gamma \neq \emptyset\}. \quad (5.5)$$

In Figures 12, 13 and 14, we have plotted the error $\epsilon^h = \left| \|u\|_{\mathcal{E}}^2 - \|u_h\|_{\mathcal{E}}^2 \right|^{1/2}$, the condition number $\kappa_2(\mathbf{A})$, and the angle between S_{FEM}^h and S_{ENR}^h , with respect to h , for the GFEM with enrichment given in (5.4) and \mathcal{R}^h as in (5.5). It is clear that the method yields the optimal order of convergence of $O(h)$, $\kappa_2(\mathbf{A}) = O(h^{-2})$, and the angle between S_{FEM}^h and S_{ENR}^h is bounded away from 0 for all the values of h considered in the experiment. Condition (4.2) is thus satisfied. We have also checked that (4.3) is satisfied (see dashed lines in Figure 13). The GFEM with enrichment given

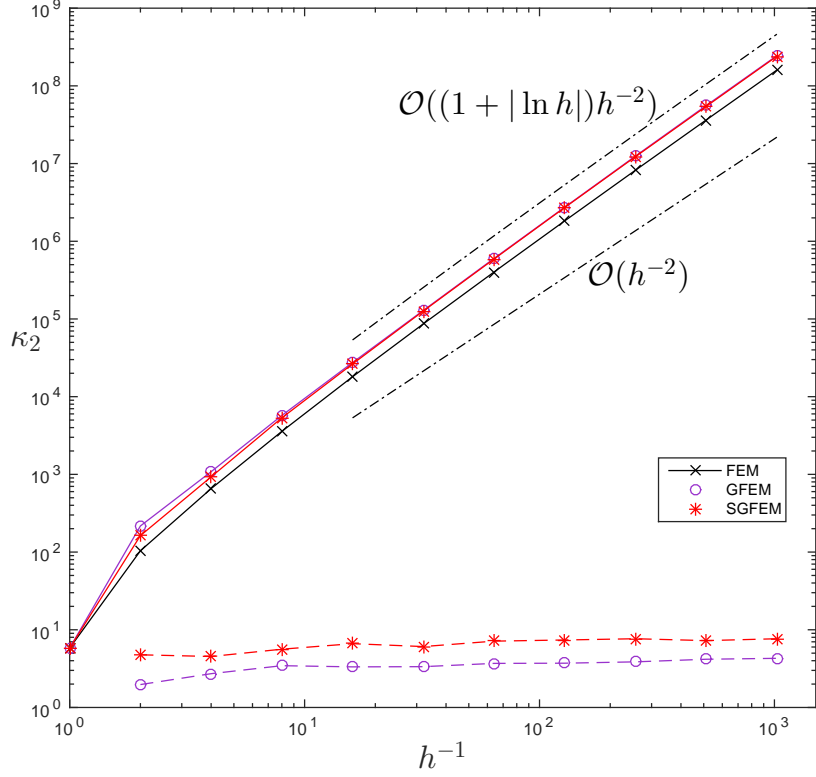


Figure 13: Condition number of the scaled stiffness matrices \mathbf{A} (solid lines) and \mathbf{A}_{22} (dashed lines) against h .

in (5.4)–(5.5) is thus an SGFEM. Once again, the angle between S_{FEM}^h and S_{ENR}^h for the SGFEM is larger than that for the GFEM. This feature is central to the iterative solution of (3.6), which we will illustrate in Section 6. We thus conclude that the SGFEM once again has all the desired features: it yields accurate approximation, it is well-conditioned, and it is robust.

We mention that if the closed curved interface Γ has a straight part, the angle between the spaces S_{FEM}^h and S_{ENR}^h associated with the GFEM considered in this section will become small when the relative distance between the element edges and the “straight part of Γ ” is small, similar to what we observed in Figure 8. This phenomenon will give rise to a much larger value of $\kappa_2(\mathbf{A})$ similar to the situation shown in Figure 7. The GFEM will thus not be stable. The SGFEM will nevertheless be stable for such interface problems. However, since in this section we considered a circular interface problem, which does not have any straight parts, the GFEM did not exhibit a behavior similar to what was observed in Figure 7 or Figure 8.

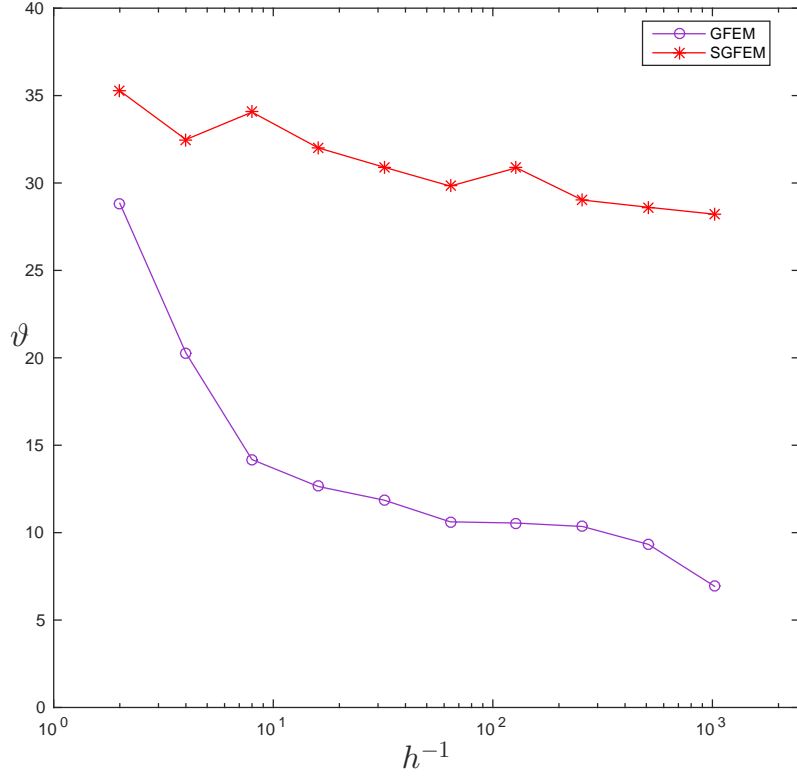


Figure 14: Angle (in degrees) between the spaces S_{FEM}^h and S_{ENR}^h against h .

6 Iterative methods

In this section, we exploit the angle condition between S_{FEM}^h and S_{ENR}^h and design appropriate iterative solvers based on previous observations. For a given error tolerance, we will compare the performance of the iterative solvers for FEM, GFEM (recall that we are only referring to M-GFEM; see Section 5) and SGFEM. The solvers designed in this section can be applied to both the straight interface problem (Section 4) and the circular interface problem (Section 5) and the conclusions we come to are very similar for the two problems. Note that the discretization error takes a different form depending on the problem: for the straight interface case, it is the classical $\epsilon^h = \|u - u_h\|_{\mathcal{E}}$, while for the circular interface case we consider instead $\epsilon^h = \left| \|u\|_{\mathcal{E}}^2 - \|u_h\|_{\mathcal{E}}^2 \right|^{1/2}$. Let $v_h \in S^h$ be an approximation of $u_h \in S^h$. *For the sake of concision, we will refer to the truncation error as $\delta := \|u_h - v_h\|_{\mathcal{E}}$ with the understanding that in the case of the circular interface problem, we actually mean $\delta = \|u_h - v_h\|_{\mathcal{E}}$.*

We start by examining the linear system associated to standard FEM (3.6) with $w \equiv 0$, yielding $S^h = S_{FEM}^h$ and $\mathbf{A} = \mathbf{A}_{11}$. The exact (discrete) solution of this system is noted u_h . We design an iterative solver based on Schur complement [12, 25]. We will denote by v_h^i the iterative solution at iteration i . There are now two sources of error: the first one is the discretization error ϵ^h due to the choice of the approximation space S^h . The second one is the truncation error due to the choice of the iterative solver, $\delta_i = \|u_h - v_h^i\|_{\mathcal{E}}$ (or $\delta_i = \|u_h - v_h^i\|_{\tilde{\mathcal{E}}}$ for the circular interface problem). The iterative solver we will use is Conjugate Gradient (CG) preconditioned by Full Multigrid (FMG). First, we recall the essential steps in multigrid methods, when they are viewed not as preconditioner but as solvers.

We start by applying a few relaxation steps (e.g., Gauss-Seidel) until the residual starts to stagnate, indicating that the low frequency content of the solution has been found. Then, we interpolate the residual onto a coarser grid and perform again a few relaxations until the residual stagnates again. This is applied until the coarsest level is reached (typically containing only a few elements), where we solve the residual equation exactly, before projecting back onto the finer grids, applying there again a few relaxations each time. This scheme is the so-called V-cycle. Applying successive V-cycles allows the truncation error δ_i to decrease geometrically from one cycle to the next. FMG is a variant of this scheme. In FMG, V-cycles are applied recursively starting on the coarsest mesh, until the finest is reached. Once again, when several FMGs are performed, the truncation error δ_i decreases geometrically, with a higher reduction factor than for a single V-cycle. However, the computational cost of an FMG is slightly larger than that of a V-cycle [12].

As stated before, we will not use FMG as a solver in this paper, but rather as a preconditioner for CG. One of the reasons for this choice is that the mathematical theory for multigrid methods is available for smooth coefficients, but it is not yet well-developed for discontinuous coefficients. As a result, we have observed that CG preconditioned by a multigrid method such as V-cycle or FMG was more robust than the multigrid method alone used as a solver. We have also noted that FMG has better computational properties than the regular V-cycle (reduction of the truncation error for the same computational work). As a result, our solver for S_{FEM}^h is CG preconditioned by FMG.

We now have to design an appropriate stopping criterion in order to decide when to stop the CG iterations. By a priori error estimation, we know that the discretization error ϵ^h behaves in some Ch^p , where p is known by the underlying properties of the PDE, the choice of the partition of unity and the choice of the enrichment. A wise stopping criterion would be to stop the iterations as soon as the truncation error δ_i becomes significantly smaller than the discretization error ϵ^h . We will show in Appendix C how to attain this by using a controlling factor. Performing more iterations will not result in a sensible improvement as the quality of the iterative solution v_h^i will mostly be driven by the discretization and not the truncation part of the error.

Due to the use of the CG solver, the truncation error is no longer expected to decay geometrically with the number of iterations. However, the use of FMG as a preconditioner reduces the “effective”

condition number $\kappa_2(\mathbf{A}_{11})$ from $O(h^{-2})$ to $O(h^{-1})$, see [27]. We thus rely on an error estimator for the truncation error δ_i based on the residual of the linear system and the “effective” spectral radius of the matrix \mathbf{A}_{11}^{-1} from (3.6) using a so-called inverse estimate. This allows us to estimate the truncation error at step i as follows

$$e^i = \frac{\|\mathbf{f}_1 - \mathbf{A}_{11}\mathbf{x}^i\|_{l^2}}{h}, \quad (6.1)$$

where the vector $\mathbf{x}^i = \mathbf{D}^{-1}\hat{\mathbf{c}}^i$ is associated to v_h^i in the approximation space S_{FEM}^h , i.e., $v_h^i = \sum_{k \in \mathcal{N}_d^h} c_k^i N_k$, and \mathbf{D} is the scaling matrix associated to \mathbf{A} from (3.6). The full derivation of this estimator can be found in Appendix C.

The algorithm developed for FEM schematically takes the form of Algorithm 1.

Data: $h, \mathbf{A}_{11}, \mathbf{f}_1, k$
Result: $v_h^{i^*}, i^*$
 $\epsilon = h^{1/2}, v_h^0 = 0, e^0 = \infty, i = 0;$
while $e^i \geq \epsilon/k$ **do**
 $i \leftarrow i + 1;$
 Compute v_h^i using initialization $v_h^{i-1};$
 Compute error estimator e^i using (6.1);
end
 $i^* = i.$

Algorithm 1: Algorithm for FEM.

Note that, as always, we work on the scaled system (3.6). The algorithm stops iterating as soon as the estimated truncation error e^i becomes significantly lower than the a priori estimated discretization error ϵ . The factor k is used to control the different constants of proportionality appearing in the intermediate calculations (see Appendix C).

Let us now consider the case of GFEM & SGFEM. Since we want to exploit the angle condition between S_{FEM}^h and S_{ENR}^h , the idea is to apply a block Gauss-Seidel iterative scheme between S_{FEM}^h and S_{ENR}^h . This defines our “outer iterations”. The system to be solved is (3.6), and by partitioning the solution \mathbf{x} in $\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2] = [x_{1,i}, x_{2,k}]_{i \in \mathcal{N}_d^h; k \in \mathcal{R}^h}$, corresponding to $u_h \in S^h = S_{FEM}^h \oplus S_{ENR}^h$ where $u_h = u_{1,h} + u_{2,h}$ with $u_{1,h} \in S_{FEM}^h$ and $u_{2,h} \in S_{ENR}^h$, we obtain

$$\text{Solve } \mathbf{A}_{11}\mathbf{x}_1^i = \mathbf{f}_1 - \mathbf{A}_{12}\mathbf{x}_2^{i-1} \text{ for } \mathbf{x}_1^i, \quad (6.2)$$

$$\text{Solve } \mathbf{A}_{22}\mathbf{x}_2^i = \mathbf{f}_2 - \mathbf{A}_{12}^T\mathbf{x}_1^i \text{ for } \mathbf{x}_2^i, \quad (6.3)$$

where \mathbf{x}^i denote successive iterates $\mathbf{x}^i = [\mathbf{x}_1^i, \mathbf{x}_2^i]$ corresponding to iterates in S^h , $v_h^i = v_{1,h}^i + v_{2,h}^i$ with $v_{1,h}^i \in S_{FEM}^h$ and $v_{2,h}^i \in S_{ENR}^h$. The truncation error $\delta_i = \|u_h - v_h^i\|_{\mathcal{E}}$ decreases geometrically with a ratio related to the angle between the spaces S_{FEM}^h and S_{ENR}^h . In fact, if q denotes this ratio, we have $q = \cos^2(\vartheta(S_{FEM}^h, S_{ENR}^h))$, see Figure 15 where $\vartheta = \pi/6$ and the truncation error is

divided at each iteration by a factor $q^{-1} = 4/3$, as are the quantities $\|v_{1,h}^i - v_{1,h}^{i-1}\|_{\mathcal{E}}$, $\|v_{2,h}^i - v_{2,h}^{i-1}\|_{\mathcal{E}}$ and $\|v_h^i - v_h^{i-1}\|_{\mathcal{E}}$.

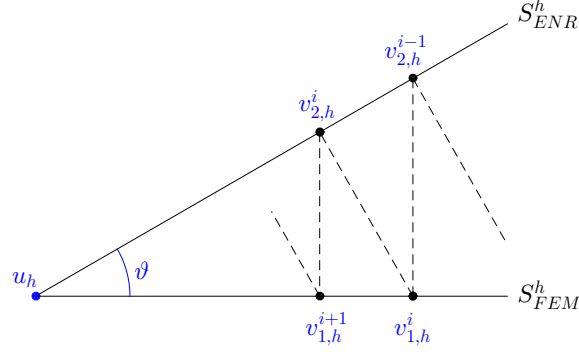


Figure 15: Outside iterations: block Gauss-Seidel scheme between S_{FEM}^h and S_{ENR}^h (6.2)–(6.3).

We can estimate the truncation error in the outer iteration $\delta_i = \|u_h - v_h^i\|_{\mathcal{E}}$ using Richardson extrapolation on the last three iterates as follows

$$e^i = \frac{1}{\frac{1}{\|v_h^i - v_h^{i-1}\|_{\mathcal{E}}} - \frac{1}{\|v_h^{i-1} - v_h^{i-2}\|_{\mathcal{E}}}}. \quad (6.4)$$

The derivation of this estimator can be found in Appendix C.

At a given stage i in the outer iteration, we now have to find $v_{1,h}^i$ using (6.2), and $v_{2,h}^i$ using (6.3). Each of these will also be done in an iterative manner, which defines our “inner iterations”.

For $v_{1,h}^i$, i.e., (6.2), we have to invert \mathbf{A}_{11} , we thus use the same solver as in FEM, that is CG preconditioned by FMG, building a sequence $w_{1,h}^{i,j}$ governed by iteration number j . The stopping criteria need not be as demanding as in FEM because we are not interested in $v_{1,h}^i$ but in $u_{1,h}$ instead. We thus design the stopping criterion so that the truncation error in the inner iteration in S_{FEM}^h , $\delta_1^{i,j} = \|v_{1,h}^i - w_{1,h}^{i,j}\|_{\mathcal{E}}$ is only a small fraction of the truncation error in S^h , $\delta_i = \|u_h - v_h^i\|_{\mathcal{E}}$.

This yields the following error estimator for $\delta_1^{i,j} = \|v_{1,h}^i - w_{1,h}^{i,j}\|_{\mathcal{E}}$

$$e_1^{i,j} = \frac{\|\mathbf{f}_1 - \mathbf{A}_{12}\mathbf{x}_2^{i-1} - \mathbf{A}_{11}\mathbf{x}_1^{i,j}\|_{l^2}}{h}. \quad (6.5)$$

The derivation of this estimator essentially follows the same as for (6.1), having replaced the right hand side \mathbf{f}_1 with $\mathbf{f}_1 - \mathbf{A}_{12}\mathbf{x}_2^{i-1}$. Once the stopping criterion is reached, after say j' iterations, we set $w_{1,h}^i := w_{1,h}^{i,j'}$ and update $w_h^i = w_{1,h}^i + w_{2,h}^i$ (in practice $w_{2,h}^i$ is to be computed in the second block of the Gauss-Seidel scheme). However, v_h^i is not available in practice, δ_i is therefore not estimated using (6.4) but rather

$$e^i = \frac{1}{\frac{1}{\|w_h^i - w_h^{i-1}\|_{\mathcal{E}}} - \frac{1}{\|w_h^{i-1} - w_h^{i-2}\|_{\mathcal{E}}}}, \quad (6.6)$$

which is the same as (6.4), having replaced v_h^i with w_h^i .

For $v_{2,h}^i$, i.e., (6.3) we have to invert \mathbf{A}_{22} , and based on condition (4.3), its condition number is bounded, so we choose to solve the equation in S_{ENR}^h using CG. We denote the successive iterates $w_{2,h}^{i,j}$. Once again, we design the stopping criterion so that the truncation error in the inner iteration in S_{ENR}^h , $\delta_2^{i,j} = \|v_{2,h}^i - w_{2,h}^{i,j}\|_{\mathcal{E}}$ is only a small fraction of the truncation error in S^h , $\delta_i = \|u_h - v_h^i\|_{\mathcal{E}}$. To estimate $\|v_{2,h}^i - w_{2,h}^{i,j}\|_{\mathcal{E}}$, we use the norm of the residual. This yields the following error estimator for $\delta_2^{i,j} = \|v_{2,h}^i - w_{2,h}^{i,j}\|_{\mathcal{E}}$

$$e_2^{i,j} = \|\mathbf{f}_2 - \mathbf{A}_{12}^T \mathbf{x}_1^i - \mathbf{A}_{22} \mathbf{x}_2^{i,j}\|_{l^2}, \quad (6.7)$$

where this time there is no factor h because the condition number of \mathbf{A}_{22} is bounded. Apart from this modification, the derivation of this estimator essentially follows the same pattern as for (6.1). Once the stopping criterion is reached, after say j'' iterations, we set $w_{2,h}^i := w_{2,h}^{i,j''}$ and update $w_h^i = w_{1,h}^i + w_{2,h}^i$. Again, note that v_h^i is not available in practice, so δ_i is estimated as (6.6) instead of (6.4).

The algorithm developed for GFEM & SGFEM schematically takes the form of Algorithm 2.

Note that, as always, we work on the scaled system (3.6). The algorithm stops the outer iterations as soon as the estimated truncation error e^i becomes significantly smaller than the a priori estimated discretization error ϵ . The factor k is there to control the different constants of proportionality appearing in the intermediate calculations (see Appendix C). During each outer iteration, the algorithm stops the inner iterations as soon as the estimated truncation error $e^{i,j}$ becomes significantly lower than the current estimated truncation error e^i . Again, the factor k' is there to control the different constants of proportionality appearing in the intermediate calculations (see Appendix C).

We show in Tables 1–3 the performances of the iterative solver described above in terms of number of iterations and of computing time when h varies. The FEM has been included in these

Data: $h, \mathbf{A}, \mathbf{f}, k, k'$
Result: $w_h^{i^*}, i^*, j^*, j'^*$
 $\epsilon = h, w_{1,h}^0 = w_{2,h}^0 = 0, e^0 = \infty, i = j^* = j'^* = 0;$
while $e^i \geq \epsilon/k$ **do**
 $i \leftarrow i + 1;$
 $w_{1,h}^{i,0} = w_{1,h}^{i-1}, e_1^{i,0} = \infty, j = 0;$
 while $e_1^{i,j} \geq e^i/k'$ **do**
 $j \leftarrow j + 1;$
 Compute $w_{1,h}^{i,j}$ using initialization $w_{1,h}^{i,j-1};$
 Compute error estimator $e_1^{i,j}$ using (6.5);
 end
 $j^* \leftarrow j^* + j, w_{1,h}^i = w_{1,h}^{i,j};$
 $w_{2,h}^{i,0} = w_{2,h}^{i-1}, e_2^{i,0} = \infty, j = 0;$
 while $e_2^{i,j} \geq e^i/k'$ **do**
 $j \leftarrow j + 1;$
 Compute $w_{2,h}^{i,j}$ using initialization $w_{2,h}^{i,j-1};$
 Compute error estimator $e_2^{i,j}$ using (6.7);
 end
 $j'^* \leftarrow j'^* + j, w_{2,h}^i = w_{2,h}^{i,j};$
 $w_h^i = w_{1,h}^i + w_{2,h}^i;$
 Compute error estimator e^i using (6.6);
end
 $i^* = i.$

Algorithm 2: Algorithm for GFEM & SGFEM.

tables only for the purpose of comparison. For a given h , the iterative solver stops when the truncation error tolerance $e^i = \frac{\epsilon}{k}$ with $k = 100$ is reached. We used $\epsilon = h^{1/2}$ for the FEM (Algorithm 1) and $\epsilon = h$ for the GFEM/SGFEM (Algorithm 2); these values are based on the a priori discretization error estimates of FEM and GFEM/SGFEM.

For GFEM & SGFEM, the number of iterations is displayed under the form $i^* (j^*, j'^*)$, where i^* is the number of outer iterations, j^* is the cumulated number of CG preconditioned by FMG iterations in S_{FEM}^h and j'^* is the cumulated number of CG iterations in S_{ENR}^h . We used a computer with 64 bits architecture, a 3.6GHz processor and 16 GB of RAM. On this computer, we measured that Matlab was able to “count” up to 150×10^6 in 1 second. The parameters for these simulations were $a_0 = 1$, $a_1 = 10$, $k = 100$ and $k' = 4$. We also mention that the relaxation scheme used in FMG was the Gauss-Seidel method and the associated stopping criterion was to stop relaxing when the l^2 -norm of the residual was higher than half of the previous one. Finally, the coarsest level in the FMG scheme is associated with the uniform mesh \mathcal{T}_h with $h = 1$, and subsequent levels are given by halving the mesh size h , and all initializations are done with the zero vector.

Table 1 illustrates the case of the straight interface with $\theta_0 = \pi/6$ and $d_0 = 1 - 1/\sqrt{2}$.

Table 1: Results of the CG preconditioned by FMG iterative solver for FEM/GFEM/SGFEM on the straight interface problem with $\theta_0 = \pi/6$.

1/h	FEM		GFEM		SGFEM	
	# it.	t (s)	# it.	t (s)	# it.	t (s)
2	2	0.0009	38 (38,59)	0.0233	8 (8,8)	0.0047
4	2	0.0028	46 (46,70)	0.0566	12 (12,12)	0.0174
8	2	0.0039	53 (53,81)	0.1109	12 (12,12)	0.0301
16	2	0.0078	59 (59,89)	0.2035	13 (13,13)	0.0449
32	3	0.0151	67 (67,103)	0.3637	13 (13,13)	0.0710
64	3	0.0268	77 (77,118)	0.7558	14 (14,14)	0.1393
128	3	0.0609	83 (113,128)	2.614	14 (14,14)	0.3374
256	4	0.2618	90 (141,139)	11.12	14 (14,14)	1.130
512	4	1.119	97 (167,150)	54.57	15 (28,15)	9.048
1024	5	5.509	103 (176,159)	227.6	15 (28,15)	35.91
2048	6	30.60	112 (229,173)	1302	15 (28,15)	158.6

Table 2 illustrates the case of the straight interface with $\theta_0 = \pi/4$ (pathological case where the interface is parallel to some of the mesh edges) and d_0 such that the relative distance between the mesh and the interface was only 10^{-3} .

Table 2: Results of the CG preconditioned by FMG iterative solver for FEM/GFEM/SGFEM on the straight interface problem with $\theta_0 = \pi/4$ and relative distance to the mesh 10^{-3} .

1/h	FEM		GFEM		SGFEM	
	# it.	t (s)	# it.	t (s)	# it.	t (s)
2	3	0.0017	44 (54,69)	0.0362	13 (13,13)	0.0088
4	3	0.0039	33 (34,49)	0.0432	17 (18,17)	0.0210
8	3	0.0056	5 (6,5)	0.0120	5 (6,5)	0.0119
16	3	0.0092	5 (6,5)	0.0193	5 (6,5)	0.0233
32	5	0.0241	5 (6,5)	0.0305	5 (6,5)	0.0305
64	5	0.0423	16 (21,23)	0.1976	7 (10,7)	0.0915
128	5	0.0955	128 (166,199)	3.772	7 (10,7)	0.2131
256	6	0.3767	119 (170,185)	13.13	9 (12,9)	0.9202
512	6	1.565	84 (133,130)	42.04	11 (14,11)	4.334
1024	7	7.609	53 (90,69)	126.1	18 (22,18)	30.62
2048	9	42.91	140 (571,218)	3128	7 (15,7)	74.25

Table 3 illustrates the case of the circular interface with $r_c = 1/\sqrt{10}$ and $(x_c, y_c) = (1/\sqrt{5}, 1/\sqrt{3})$.

The # it. and t(s) for FEM for a given h in Tables 1–3 are much less than for GFEM/SGFEM. This is expected as the error in FEM for a given h ($O(h^{1/2})$) is much greater than that of GFEM/SGFEM ($O(h)$). Moreover, we observe that the computational time t scales a little over quadratically with the mesh size and we have roughly $t = O(h^{-2.3})$ (except for the pathological case of the straight interface problem with $\theta_0 = \pi/4$ and “small” relative distance to the mesh). This rate is slightly over the optimal rate of $O(h^{-2})$ because although we use efficient solvers, our stopping criteria are somewhat pessimistic since they rely on the norm of the residual. Concerning the pathological case of the interface parallel to the mesh edges, Table 2 reveals that GFEM is not

Table 3: Results of the CG preconditioned by FMG iterative solver for FEM/GFEM/SGFEM on the circular interface problem.

1/h	FEM		GFEM		SGFEM	
	# it.	t (s)	# it.	t (s)	# it.	t (s)
2	2	0.0009	6 (6,6)	0.0036	6 (6,6)	0.0035
4	2	0.0021	22 (22,22)	0.0255	8 (8,8)	0.0093
8	3	0.0057	23 (23,23)	0.0478	10 (10,10)	0.0200
16	3	0.0092	49 (49,64)	0.1681	10 (10,10)	0.0326
32	4	0.0192	62 (63,94)	0.3506	13 (13,13)	0.0667
64	4	0.0344	83 (115,131)	1.109	16 (17,16)	0.1610
128	5	0.0970	92 (131,147)	2.940	15 (18,15)	0.3932
256	6	0.3669	100 (165,166)	11.25	16 (30,16)	1.995
512	6	1.445	120 (271,201)	74.48	17 (32,17)	9.169
1024	7	7.217	142 (338,244)	398.2	20 (40,20)	47.01
2048	8	37.41	156 (387,270)	1976	22 (51,22)	259.4

stable, while SGFEM is. Remember that in Figures 7 and 8 we saw that the condition number of (M-)GFEM blew up as the interface was getting closer and closer to the mesh edges, while the angle was going to 0. Conversely, SGFEM was stable in this situation.

We also note that for the three situations shown in Tables 1-3, many more outer iterations are needed for GFEM than for SGFEM. This is a direct consequence of the fact that the angle for SGFEM is larger than that for the GFEM, as indicated in Figures 6 and 14. Indeed, if we were to solve the block Gauss-Seidel system (6.2)–(6.3) exactly (i.e., we only had outer iterations) and if there had been only one angle between the spaces (i.e., the smallest) then, for the same decrease in the truncation error, for each outer iteration in SGFEM we would have needed $n = \frac{\log(\cos \vartheta_2)}{\log(\cos \vartheta_1)}$ outer iterations in GFEM, where ϑ_1 denotes the angle for GFEM and ϑ_2 denotes the angle for SGFEM. Indeed, each iteration will result in a reduction of the truncation error by a factor $q_1 = \cos^2 \vartheta_1$ for GFEM and $q_2 = \cos^2 \vartheta_2$ for SGFEM. As a result, n iterations for GFEM will decrease the error by q_1^n . Solving $q_1^n = q_2$ for n yields the above ratio in terms of angles ϑ_1, ϑ_2 . We have in fact solved the block Gauss-Seidel system (6.2)–(6.3) by using direct solvers for the blocks and the stopping criterion (6.4). We did observe the role of the angle, i.e., SGFEM yielded a speed-up of roughly 6-9 times compared to the GFEM. However our iterative solver with two estimators performed much faster (by a factor of about 7 for the values of h considered in our experiments) on both GFEM and SGFEM than solving the block Gauss-Seidel system (6.2)–(6.3) directly. We do not show the results as in this paper we are only concerned with our iterative scheme. Note that our error estimators estimate the truncation error as well as the discretization error and then “balance” these two errors; these estimators are very fast to compute.

However, since we do not solve (6.2)–(6.3) exactly and use estimates to decide when to stop the inner iterations, this ratio n mentioned in the last paragraph is slightly perturbed. Recall also that the angle we discuss in this paper is in fact the smallest between the spaces. There are other, larger, angles that could affect the estimation of the truncation error since we use Richardson extrapolation

to estimate δ_i using (6.6). Since the angle for SGFEM is larger than for GFEM, this “pollution” by larger angles affects GFEM more than SGFEM. As a result, the risk of under-resolving system (3.6) is higher for GFEM than SGFEM.

We also considered other solvers not shown here (e.g., V-cycles as a solver, FMG followed by V-cycles as a solver, CG preconditioned by some V-cycles), however, the presented solver – CG preconditioned by FMG – was found to be the most robust and computationally efficient. The reason is that CG preconditioned by multigrid is more robust than multigrid alone, and that FMG is more efficient than V-cycles.

Let us verify if the iterative solutions have converged to discretization accuracy. We now denote by $v_h \in S^h$ the iterated solution of (3.6) yielded by the chosen iterative solver and by $\hat{\epsilon}^h = \|u - v_h\|_{\mathcal{E}}$ the total error due to both discretization and truncation. Thanks to Galerkin orthogonality between $u - u_h$ and $u_h - v_h$ in the $B(\cdot, \cdot)$ inner product, it holds, for the straight interface problem

$$\begin{aligned} (\hat{\epsilon}^h)^2 &= \|u - v_h\|_{\mathcal{E}}^2, \\ &= \|u - u_h\|_{\mathcal{E}}^2 + \|u_h - v_h\|_{\mathcal{E}}^2, \\ &= (\epsilon^h)^2 + \delta^2, \end{aligned}$$

so that the total error $\hat{\epsilon}^h$ can be orthogonally decomposed into discretization error ϵ^h and truncation error δ . Similarly, for the circular interface problem we set $(\hat{\epsilon}^h)^2 = (\epsilon^h)^2 + \delta^2$, so that

$$(\hat{\epsilon}^h)^2 = \left| \|u\|_{\mathcal{E}}^2 - \|u_h\|_{\mathcal{E}}^2 \right| + \|u_h - v_h\|_{\mathcal{E}}^2.$$

In practice, we also solve (3.6) using a direct solver in order to have u_h . We can thus compute ϵ^h as well as $\hat{\epsilon}^h$.

Figure 16 displays the evolution of the (relative) total error $\hat{\epsilon}^h$ as the computational time t varies, on the circular interface problem (same parameters as for Table 3). We observe that the error $\hat{\epsilon}^h$ scales with the computational time t for FEM as $\hat{\epsilon}^h = 0.04 \times t^{-0.22}$. For GFEM, we have roughly $\hat{\epsilon}^h = 0.015 \times t^{-0.43}$. For SGFEM, we have $\hat{\epsilon}^h = 0.006 \times t^{-0.43}$. The optimal rates would be $O(t^{-1/4})$ for FEM and $O(t^{-1/2})$ for GFEM & SGFEM. Moreover, the SGFEM is roughly eight or nine times faster than the GFEM for the same accuracy. As discussed earlier, this is a direct consequence of the angle property. We also mention that for the straight interface problem (with the same parameters as for Table 1) for FEM we have $\hat{\epsilon}^h = 0.0275 \times t^{-0.22}$. For GFEM, we have roughly $\hat{\epsilon}^h = 0.0075 \times t^{-0.43}$. For SGFEM, we have $\hat{\epsilon}^h = 0.0035 \times t^{-0.43}$. This in turn tells us that the SGFEM is roughly six times faster than the GFEM for the same accuracy on this problem.

To verify if our iterative solutions have indeed converged to discretization accuracy, we show in Figure 17 the evolution of the ratio iterative solution error over discretization error $i^h := \frac{\hat{\epsilon}^h}{\epsilon^h} = \sqrt{1 + \left(\frac{\delta}{\epsilon^h}\right)^2}$ (on the circular interface problem with the same parameters as before). We observe

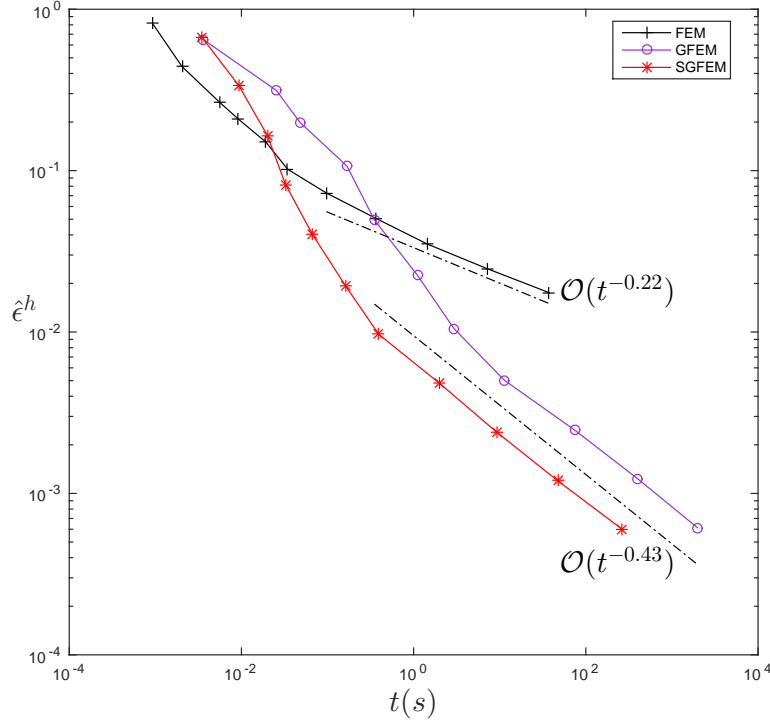


Figure 16: Evolution of the (relative) error $\hat{\epsilon}^h$ as the computational time varies on the circular interface problem.

that the ratio i^h is very close to unity for all considered methods. There is no appreciable difference for FEM or SGFEM, while for GFEM the ratio i^h is less than 2% over unity. Our solutions have thus converged to discretization accuracy. As discussed before, GFEM performs worse than SGFEM because of the angle property and the pollution by larger angles in the estimation of the truncation error.

A last remark has to be made about the solutions yielded by the iterative solvers in the case where the exact solution $u \in \mathcal{E}$ and the discrete solution $u_h \in S^h$ are unknown. Once the iterative solutions have been computed they can be used to design an error estimator. Indeed, neither the solver nor the stopping criteria rely on the knowledge of u or u_h . Using Richardson extrapolation this time on the norm of v_h allows us to extrapolate when $h \rightarrow 0$ to obtain some η approximating $\|u\|_{\mathcal{E}}$. Indeed

$$\begin{aligned} \|u\|_{\mathcal{E}} &= \|v_h\|_{\mathcal{E}} + \varepsilon_h, \\ \|u\|_{\mathcal{E}} &= \|v_{2h}\|_{\mathcal{E}} + \varepsilon_{2h}. \end{aligned}$$

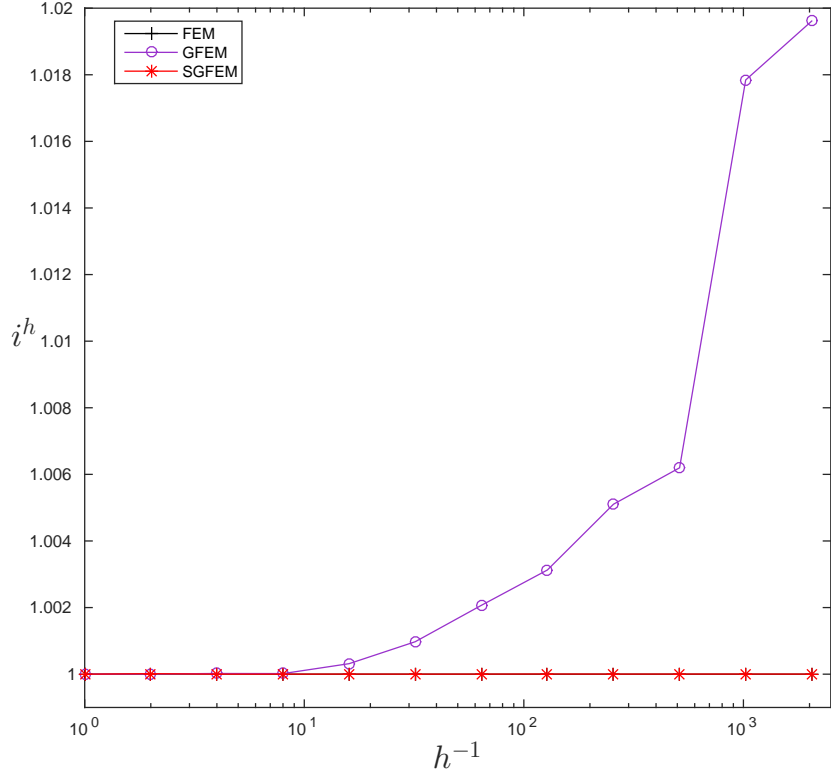


Figure 17: Evolution of the ratio iterative solution error over discretization error against h .

Next, assuming $\varepsilon_h = Ch^p$ where p is known by the underlying properties of the PDE, the choice of the partition of unity and the choice of the enrichment, we have

$$C = \frac{\|v_h\|_{\mathcal{E}} - \|v_{2h}\|_{\mathcal{E}}}{(2^p - 1)h^p}.$$

As a result, we can estimate $\|u\|_{\mathcal{E}}$ as

$$\eta = \|v_h\|_{\mathcal{E}} + Ch^p.$$

Finally, by assuming Galerkin orthogonality, we have an error estimator using $\bar{\epsilon}^h := |\eta^2 - \|v_h\|_{\mathcal{E}}^2|^{1/2}$. The graph of $\bar{\epsilon}^h$ as the computational time t varies is similar to Figure 16 and we do not include it here. The efficiency of this estimator can then be assessed by computing $\hat{i}^h := \frac{\bar{\epsilon}^h}{\epsilon^h}$. The results are presented in Figure 18 (on the circular interface problem with the same parameters as before). We observe that the efficiency of the error estimator \hat{i}^h stays close to unity for all considered methods:

\hat{i}^h is within the 3% range for FEM and SGFEM, 40% for GFEM. Once again, GFEM performs worse than SGFEM because of the angle property and the pollution by larger angles in the estimation of the truncation error.

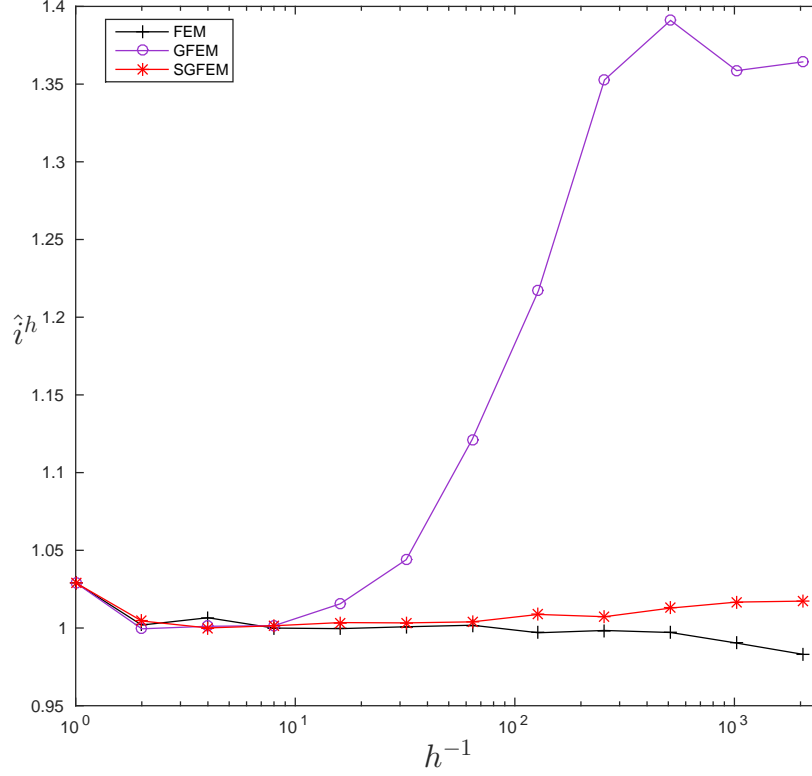


Figure 18: Efficiency of the error estimator i^h as h varies.

We summarize the results mentioned above about the differences between GFEM and SGFEM in the following table, where for a given relative error tolerance $\tau = 5\%, 1\%, 0.1\%$, we show the computed relative total error $\hat{\epsilon}^h$, the computed relative discretization error ϵ^h , the computed relative truncation error δ , the computed relative extrapolated error $\bar{\epsilon}^h$, time $t(s)$ and the efficiency indexes i^h and \hat{i}^h , on the circular interface problem with the same parameters as before.

τ	GFEM						
	$\hat{\epsilon}^h$	ϵ^h	δ	$\bar{\epsilon}^h$	$t(s)$	i^h	\hat{i}^h
5%	4.95%	4.94%	0.22%	5.16%	0.3506	1.001	1.044
1%	1.05%	1.05%	0.083%	1.28%	2.940	1.003	1.217
0.1%	0.123%	0.121%	0.023%	0.168%	398.2	1.018	1.359

	SGFEM						
τ	$\hat{\epsilon}^h$	ϵ^h	δ	$\bar{\epsilon}^h$	$t(s)$	i^h	\hat{i}^h
5%	4.03%	4.03%	0.0069%	4.05%	0.0667	1.000	1.003
1%	0.967%	0.967%	0.0032%	0.976%	0.3932	1.000	1.009
0.1%	0.120%	0.120%	0.00032%	0.122%	47.01	1.000	1.017

Conclusion

In this paper, we have considered several GFEMs on an interface problem. The proposed GFEMs differed with regards to the choice of enrichment function and space. We have illustrated how these different choices could allow recovering the optimal order of convergence of the underlying partition of unity. However, as we have shown, this so-called approximation property is not the only important feature of a well-designed GFEM, and it is equally important that the GFEM be well-conditioned in order to be able to solve the resulting linear system efficiently. To this end, we have highlighted the importance of the so-called angle between the approximation spaces. We emphasized how this angle was related to the conditioning of GFEM and to the stability of the method with respect to varying criteria, such as the position of the interface with respect to the mesh. Finally, we have showed in the last section that this angle property could be exploited within a block Gauss-Seidel iterative scheme between the approximation spaces, resulting in large savings in computational times.

A Angle between subspaces

In this first annex, we derive the formula yielding the angle ϑ between subspaces S_1 and S_2 , of dimensions m and n respectively, in the sense of the inner product $B(\cdot, \cdot)$.

Let $S = S_1 + S_2 = \{s = (s_1, s_2) : s_1 \in S_1, s_2 \in S_2\}$ and P_1 be the orthogonal projection operator on S_1 . Then the angle ϑ is defined by

$$\cos \vartheta = \max_{s_2 \in S_2} \frac{\|P_1(s_2)\|_{\mathcal{E}}}{\|s_2\|_{\mathcal{E}}}.$$

By definition of the projection P_1 , we have

$$B(P_1(s_2), s_1) = B(s_2, s_1), \quad \forall s_1 \in S_1. \quad (\text{A.1})$$

Identifying $s_1 \approx X_1 \in \mathbb{R}^m$, $s_2 \approx X_2 \in \mathbb{R}^n$ and $P_1(s_2) \approx W \in \mathbb{R}^m$, we can then use the matrix \mathbf{A} defined in (3.6), so problem (A.1) reads

$$(W^T, 0)\mathbf{A}(X_1, 0)^T = (0, X_2^T)\mathbf{A}(X_1, 0)^T, \quad \forall X_1 \in \mathbb{R}^m.$$

With the sub-matrices defined in (3.6), we have

$$W^T \mathbf{A}_{11} X_1 = X_2^T \mathbf{A}_{21} X_1, \quad \forall X_1 \in \mathbb{R}^m$$

which leads to (\mathbf{A} being symmetric)

$$W = \mathbf{A}_{11}^{-1} \mathbf{A}_{12} X_2.$$

Then, one can simply derive

$$\begin{aligned} \|P_1(s_2)\|_{\mathcal{E}}^2 &= W^T \mathbf{A}_{11} W, \\ &= X_2^T \mathbf{A}_{21} \mathbf{A}_{11}^{-1} \mathbf{A}_{12} X_2, \end{aligned}$$

and we obviously have

$$\|s_2\|_{\mathcal{E}}^2 = X_2^T \mathbf{A}_{22} X_2.$$

Hence

$$\cos^2 \vartheta = \max_{X_2 \in \mathbb{R}^n} \frac{X_2^T \mathbf{A}_{21} \mathbf{A}_{11}^{-1} \mathbf{A}_{12} X_2}{X_2^T \mathbf{A}_{22} X_2}, \quad (\text{A.2})$$

which leads to the generalized eigenvalue problem

$$\mathbf{A}_{21} \mathbf{A}_{11}^{-1} \mathbf{A}_{12} = \lambda \mathbf{A}_{22}.$$

The largest eigenvalue of this problem is equal to $\max_{X_2 \in \mathbb{R}^n} \frac{X_2^T \mathbf{A}_{21} \mathbf{A}_{11}^{-1} \mathbf{A}_{12} X_2}{X_2^T \mathbf{A}_{22} X_2}$, so one can then find ϑ using (A.2).

B Error induced by the perturbation

In this second annex, we show

$$\left| \|u - u_h\|_{\mathcal{E}}^2 - \|u\|_{\mathcal{E}}^2 - \|u_h\|_{\mathcal{E}}^2 \right| \leq O(h^k),$$

where $k = 3/2$ for FEM and $k = 2$ for GFEM & SGFEM, $u \in \mathcal{E}$ is the (exact) solution of (3.1), $u_h \in S^h$ is the (discrete perturbed) solution of (5.3), $\|\cdot\|_{\mathcal{E}} := B(\cdot, \cdot)^{1/2}$ is the usual energy norm and $\|\cdot\|_{\tilde{\mathcal{E}}} := \tilde{B}(\cdot, \cdot)^{1/2}$ is the perturbed energy norm. We also introduce $\tilde{u} \in \mathcal{E}$ the (exact perturbed) solution of (5.2).

The proof is divided into two parts: in the first part, we show that $\left| \|u\|_{\mathcal{E}}^2 - \|\tilde{u}\|_{\tilde{\mathcal{E}}}^2 \right| = O(h^2)$. In the second part, we use this result to show that $\left| \|u - u_h\|_{\mathcal{E}}^2 - \|u\|_{\mathcal{E}}^2 - \|u_h\|_{\mathcal{E}}^2 \right| \leq O(h^k)$, where $k = 3/2$ for FEM and $k = 2$ for GFEM & SGFEM.

We have two weak formulations: one for the original problem (3.1) and one for the perturbed problem (5.2)

$$B(u, v) = \int_{\Omega} a \nabla u \cdot \nabla v \, d\mathbf{x} = \int_{\partial\Omega} g_N v \, ds, \quad \forall v \in \mathcal{E}, \quad (\text{B.1})$$

$$\tilde{B}(\tilde{u}, v) = \int_{\Omega} \tilde{a} \nabla \tilde{u} \cdot \nabla v \, d\mathbf{x} = \int_{\partial\Omega} g_N v \, ds, \quad \forall v \in \mathcal{E}. \quad (\text{B.2})$$

Let us first show that the two norms induced by $B(\cdot, \cdot)$ and $\tilde{B}(\cdot, \cdot)$ are equivalent on \mathcal{E} . For all $v \in \mathcal{E}$, we have

$$\begin{aligned} \beta_0 \int_{\Omega} \nabla v \cdot \nabla v \, d\mathbf{x} &\leq B(v, v) \leq \beta_1 \int_{\Omega} \nabla v \cdot \nabla v \, d\mathbf{x}, \\ \beta_0 |v|_{H^1}^2 &\leq B(v, v) \leq \beta_1 |v|_{H^1}^2, \end{aligned} \quad (\text{B.3})$$

where β_0, β_1 were defined in Section 2 and represent bounds on the coefficient a . As a result, the norm induced by $B(\cdot, \cdot)$ and the H^1 semi-norm are equivalent on \mathcal{E} . The same reasoning follows for $\tilde{B}(\cdot, \cdot)$, which is also equivalent to $|\cdot|_{H^1}$ on \mathcal{E} . As a result, $B(\cdot, \cdot)$ and $\tilde{B}(\cdot, \cdot)$ are equivalent on \mathcal{E} . We further emphasize that the bounds appearing in (B.3) are independent of h . It follows that $B(v, v) = O(h^p)$ holds if and only if $\tilde{B}(v, v) = O(h^p)$ holds.

Now, considering (B.1) and (B.2), we have

$$B(u, v) = \tilde{B}(\tilde{u}, v), \quad \forall v \in \mathcal{E}. \quad (\text{B.4})$$

We immediately obtain

$$B(u, u) = \tilde{B}(\tilde{u}, u), \text{ and } B(u, \tilde{u}) = \tilde{B}(\tilde{u}, \tilde{u}). \quad (\text{B.5})$$

We also have

$$\begin{aligned} B(u - \tilde{u}, v) &= B(u, v) - B(\tilde{u}, v), \text{ and, using (B.4),} \\ &= \tilde{B}(\tilde{u}, v) - B(\tilde{u}, v), \\ &= \int_{\Omega} (\tilde{a} - a) \nabla \tilde{u} \cdot \nabla v \, d\mathbf{x}, \\ &= (a_0 - a_1) \int_{\omega} \nabla \tilde{u} \cdot \nabla v \, d\mathbf{x}, \quad \forall v \in \mathcal{E}. \end{aligned} \quad (\text{B.6})$$

Using Cauchy-Schwarz inequality, it yields

$$\begin{aligned} |B(u - \tilde{u}, v)| &\leq |a_0 - a_1| \left(\int_{\omega} \nabla \tilde{u} \cdot \nabla \tilde{u} \, d\mathbf{x} \right)^{1/2} \left(\int_{\omega} \nabla v \cdot \nabla v \, d\mathbf{x} \right)^{1/2}, \\ &\leq |a_0 - a_1| \|\nabla \tilde{u}\|_{L^2(\omega)} \|\nabla v\|_{L^2(\omega)}. \end{aligned} \quad (\text{B.7})$$

From this we can obtain

$$\begin{aligned} \left| B(u, u) - \tilde{B}(\tilde{u}, \tilde{u}) \right| &= |B(u, u) - B(u, \tilde{u})|, \text{ using (B.5),} \\ &= |B(u, u - \tilde{u})|, \\ &\leq |a_0 - a_1| \|\nabla \tilde{u}\|_{L^2(\omega)} \|\nabla u\|_{L^2(\omega)}, \text{ using (B.7) with } v = u, \\ &\leq |a_0 - a_1| \mu(\omega)^{1/2} \|\nabla \tilde{u}\|_{L^2(\omega)} \|\nabla u\|_{L^\infty(\omega)}, \end{aligned} \quad (\text{B.8})$$

where we have used $\|\nabla u\|_{L^2(\omega)} \leq \mu(\omega)^{1/2} \|\nabla u\|_{L^\infty(\omega)}$ and the term on the right hand side is finite because our manufactured solution u does not have any singularity in $\overline{\Omega}$. Now, using $\mu(\omega) = O(h^2)$, we get

$$\left| B(u, u) - \tilde{B}(\tilde{u}, \tilde{u}) \right| = O(h). \quad (\text{B.9})$$

Similarly to (B.6)

$$\begin{aligned} \tilde{B}(u - \tilde{u}, v) &= \tilde{B}(u, v) - \tilde{B}(\tilde{u}, v), \text{ and, using (B.4),} \\ &= \tilde{B}(u, v) - B(u, v), \\ &= (a_0 - a_1) \int_{\omega} \nabla u \cdot \nabla v \, d\mathbf{x}, \quad \forall v \in \mathcal{E}. \end{aligned} \quad (\text{B.10})$$

Let us now look at

$$\begin{aligned} \tilde{B}(u - \tilde{u}, u - \tilde{u}) &= \tilde{B}(u - \tilde{u}, u) - \tilde{B}(u - \tilde{u}, \tilde{u}), \\ &= (a_0 - a_1) \|\nabla u\|_{L^2(\omega)}^2 - \tilde{B}(u - \tilde{u}, \tilde{u}), \text{ using (B.10) with } v = u, \\ &= (a_0 - a_1) \|\nabla u\|_{L^2(\omega)}^2 + \tilde{B}(\tilde{u}, \tilde{u}) - \tilde{B}(u, \tilde{u}), \\ &= (a_0 - a_1) \|\nabla u\|_{L^2(\omega)}^2 + \tilde{B}(\tilde{u}, \tilde{u}) - B(u, u), \text{ using (B.5).} \end{aligned}$$

Thus

$$\begin{aligned} \tilde{B}(u - \tilde{u}, u - \tilde{u}) &\leq |a_0 - a_1| \|\nabla u\|_{L^2(\omega)}^2 + \left| B(u, u) - \tilde{B}(\tilde{u}, \tilde{u}) \right|, \text{ using triangular inequality,} \\ &\leq \left| B(u, u) - \tilde{B}(\tilde{u}, \tilde{u}) \right| + O(h^2), \end{aligned} \quad (\text{B.11})$$

where again we have used $\|\nabla u\|_{L^2(\omega)} \leq \mu(\omega)^{1/2} \|\nabla u\|_{L^\infty(\omega)}$ and $\mu(\omega) = O(h^2)$.

We now have most of the ingredients to prove that $\left| B(u, u) - \tilde{B}(\tilde{u}, \tilde{u}) \right| = O(h^2)$. We will proceed by induction. Let us consider the two statements

$$\left\{ \begin{array}{l} \left| B(u, u) - \tilde{B}(\tilde{u}, \tilde{u}) \right| = O(h). \\ \left| B(u, u) - \tilde{B}(\tilde{u}, \tilde{u}) \right| = O(h^p) \Rightarrow \left| B(u, u) - \tilde{B}(\tilde{u}, \tilde{u}) \right| = O(h^{1+p/2}). \end{array} \right.$$

The first statement has already been proven in (B.9). Let us prove the second. Assume that $\left| B(u, u) - \tilde{B}(\tilde{u}, \tilde{u}) \right| = O(h^p)$ holds for some $1 \leq p \leq 2$. Then, using the induction assumption in (B.11), we have

$$\begin{aligned} \tilde{B}(u - \tilde{u}, u - \tilde{u}) &\leq O(h^p) + O(h^2), \\ &\leq O(h^p). \end{aligned}$$

Then, by the equivalence of $B(\cdot, \cdot)$ and $\tilde{B}(\cdot, \cdot)$ on \mathcal{E} , we also have

$$B(u - \tilde{u}, u - \tilde{u}) \leq O(h^p). \quad (\text{B.12})$$

Now, consider the following

$$\begin{aligned}
|B(u - \tilde{u}, u)| &= |B(u, u) - B(\tilde{u}, u)|, \\
&= \left| B(u, u) - \tilde{B}(\tilde{u}, \tilde{u}) \right|, \text{ using (B.5),} \\
&\leq O(h^p), \text{ using the induction assumption again.}
\end{aligned} \tag{B.13}$$

Now, using (B.6) with $v = \tilde{u}$

$$\begin{aligned}
|a_0 - a_1| \|\nabla \tilde{u}\|_{L^2(\omega)}^2 &= |B(u - \tilde{u}, \tilde{u})|, \\
&= |B(u - \tilde{u}, u) - B(u - \tilde{u}, u - \tilde{u})|, \\
&\leq |B(u - \tilde{u}, u)| + |B(u - \tilde{u}, u - \tilde{u})|, \text{ by the triangular inequality,} \\
&\leq O(h^p), \text{ using (B.12) and (B.13).}
\end{aligned}$$

Which yields

$$\|\nabla \tilde{u}\|_{L^2(\omega)} \leq O(h^{p/2}). \tag{B.14}$$

Recall (B.8)

$$\begin{aligned}
\left| B(u, u) - \tilde{B}(\tilde{u}, \tilde{u}) \right| &\leq |a_0 - a_1| \mu(\omega)^{1/2} \|\nabla \tilde{u}\|_{L^2(\omega)} \|\nabla u\|_{L^\infty(\omega)}, \\
&\leq O(h^{1+p/2}), \text{ using (B.14) and } \mu(\omega) = O(h^2),
\end{aligned}$$

which is the desired result. Applying it inductively starting at $p = 1$ yields $\left| B(u, u) - \tilde{B}(\tilde{u}, \tilde{u}) \right| = O(h^2)$. Equivalently, we have $\left| \|u\|_{\mathcal{E}}^2 - \|\tilde{u}\|_{\mathcal{E}}^2 \right| = O(h^2)$. In particular, all the relations written in the induction proof hold for $p = 2$.

Now, let us prove the second result, which is

$$\left| \|u - u_h\|_{\mathcal{E}}^2 - \left| \|u\|_{\mathcal{E}}^2 - \|u_h\|_{\mathcal{E}}^2 \right| \right| \leq O(h^k),$$

where $k = 3/2$ for FEM and $k = 2$ for GFEM & SGFEM and $u_h \in S^h$ is the solution of the following variational problem

$$\tilde{B}(u_h, v) = F(v), \quad \forall v \in S^h.$$

We split the result into these two inequalities for the sake of clarity: we need to show that

$$\begin{aligned}
\|u - u_h\|_{\mathcal{E}}^2 &\leq \left| \|u\|_{\mathcal{E}}^2 - \|u_h\|_{\mathcal{E}}^2 \right| + O(h^k), \text{ and,} \\
\left| \|u\|_{\mathcal{E}}^2 - \|u_h\|_{\mathcal{E}}^2 \right| &\leq \|u - u_h\|_{\mathcal{E}}^2 + O(h^k).
\end{aligned}$$

Using the triangular inequality, it holds

$$\begin{aligned}\|u - u_h\|_{\mathcal{E}} &\leq \|u - \tilde{u}\|_{\mathcal{E}} + \|\tilde{u} - u_h\|_{\mathcal{E}}, \\ &\leq \|\tilde{u} - u_h\|_{\mathcal{E}} + O(h), \text{ using (B.12) with } p = 2.\end{aligned}\tag{B.15}$$

Next, let us consider the difference

$$\begin{aligned}\left| \|\tilde{u} - u_h\|_{\mathcal{E}}^2 - \|\tilde{u} - u_h\|_{\tilde{\mathcal{E}}}^2 \right| &= \left| \int_{\Omega} (a - \tilde{a}) \nabla(\tilde{u} - u_h) \cdot \nabla(\tilde{u} - u_h) d\mathbf{x} \right|, \\ &= |a_1 - a_0| \|\nabla(\tilde{u} - u_h)\|_{L^2(\omega)}^2.\end{aligned}$$

Using (B.14) with $p = 2$ yields $\|\nabla \tilde{u}\|_{L^2(\omega)}^2 \leq O(h^2)$. And $\|\nabla u_h\|_{L^2(\omega)}^2 \leq \mu(\omega) \|\nabla u_h\|_{L^\infty(\omega)}^2 \leq O(h^2)$ as well since $u_h \in S^h$ does not have any singularity in $\bar{\Omega}$. As a result

$$\left| \|\tilde{u} - u_h\|_{\mathcal{E}}^2 - \|\tilde{u} - u_h\|_{\tilde{\mathcal{E}}}^2 \right| \leq O(h^2).\tag{B.16}$$

Then, thanks to Galerkin orthogonality between $\tilde{u} - u_h$ and u_h in the $\tilde{B}(\cdot, \cdot)$ inner-product, we have

$$\|\tilde{u} - u_h\|_{\tilde{\mathcal{E}}}^2 = \|\tilde{u}\|_{\tilde{\mathcal{E}}}^2 - \|u_h\|_{\tilde{\mathcal{E}}}^2.$$

As a result, it holds

$$\|\tilde{u} - u_h\|_{\mathcal{E}}^2 \leq \|\tilde{u}\|_{\tilde{\mathcal{E}}}^2 - \|u_h\|_{\tilde{\mathcal{E}}}^2 + O(h^2).\tag{B.17}$$

By a priori error estimation, we have

$$\|\tilde{u} - u_h\|_{\tilde{\mathcal{E}}} = O(h^p),$$

where $p = 1/2$ for FEM and $p = 1$ for GFEM & SGFEM. By equivalence of $B(\cdot, \cdot)$ and $\tilde{B}(\cdot, \cdot)$ on \mathcal{E} , it holds

$$\|\tilde{u} - u_h\|_{\mathcal{E}} = O(h^p).\tag{B.18}$$

Hence, starting with (B.15)

$$\begin{aligned}\|u - u_h\|_{\mathcal{E}}^2 &\leq (\|\tilde{u} - u_h\|_{\mathcal{E}} + O(h))^2, \\ &\leq \|\tilde{u} - u_h\|_{\mathcal{E}}^2 + \|\tilde{u} - u_h\|_{\mathcal{E}} O(h) + O(h^2), \\ &\leq \|\tilde{u} - u_h\|_{\mathcal{E}}^2 + O(h^{1+p}), \text{ using (B.18),} \\ &\leq \|\tilde{u}\|_{\tilde{\mathcal{E}}}^2 - \|u_h\|_{\tilde{\mathcal{E}}}^2 + O(h^{1+p}), \text{ using (B.17).}\end{aligned}$$

Finally, using the first part of the proof: $\left| \|u\|_{\mathcal{E}}^2 - \|\tilde{u}\|_{\mathcal{E}}^2 \right| = O(h^2)$, it holds by triangular inequality

$$\begin{aligned} \|u - u_h\|_{\mathcal{E}}^2 &\leq \left| \|\tilde{u}\|_{\mathcal{E}}^2 - \|u\|_{\mathcal{E}}^2 \right| + \left| \|u\|_{\mathcal{E}}^2 - \|u_h\|_{\mathcal{E}}^2 \right| + O(h^{1+p}), \\ &\leq \left| \|u\|_{\mathcal{E}}^2 - \|u_h\|_{\mathcal{E}}^2 \right| + O(h^{1+p}), \end{aligned}$$

which shows the first inequality since $1 + p = 3/2$ for FEM and 2 for GFEM & SGFEM.

For the second inequality, we begin with

$$\begin{aligned} \left| \|u\|_{\mathcal{E}}^2 - \|u_h\|_{\mathcal{E}}^2 \right| &= \left| \|u\|_{\mathcal{E}}^2 - \|\tilde{u}\|_{\mathcal{E}}^2 + \|\tilde{u}\|_{\mathcal{E}}^2 - \|u_h\|_{\mathcal{E}}^2 \right|, \\ &\leq \left| \|u\|_{\mathcal{E}}^2 - \|\tilde{u}\|_{\mathcal{E}}^2 \right| + \|\tilde{u}\|_{\mathcal{E}}^2 - \|u_h\|_{\mathcal{E}}^2, \text{ by triangular inequality,} \\ &\leq \|\tilde{u} - u_h\|_{\mathcal{E}}^2 + O(h^2), \end{aligned}$$

where again, we have used the first part of the proof: $\left| \|u\|_{\mathcal{E}}^2 - \|\tilde{u}\|_{\mathcal{E}}^2 \right| = O(h^2)$ and then Galerkin orthogonality between $\tilde{u} - u_h$ and u_h in the $\tilde{B}(\cdot, \cdot)$ inner-product. Now, using the same kind of reasoning as for the first inequality, we have

$$\begin{aligned} \left| \|u\|_{\mathcal{E}}^2 - \|u_h\|_{\mathcal{E}}^2 \right| &\leq \|\tilde{u} - u_h\|_{\mathcal{E}}^2 + O(h^2), \\ &\leq \|\tilde{u} - u_h\|_{\mathcal{E}}^2 + O(h^2), \text{ using (B.16),} \\ &\leq (\|u - \tilde{u}\|_{\mathcal{E}} + \|u - u_h\|_{\mathcal{E}})^2 + O(h^2), \text{ using the triangular inequality,} \\ &\leq \|u - \tilde{u}\|_{\mathcal{E}}^2 + \|u - u_h\|_{\mathcal{E}}^2 + 2\|u - \tilde{u}\|_{\mathcal{E}}\|u - u_h\|_{\mathcal{E}} + O(h^2), \\ &\leq \|u - u_h\|_{\mathcal{E}}^2 + O(h)\|u - u_h\|_{\mathcal{E}} + O(h^2), \text{ using (B.12) with } p = 2. \end{aligned}$$

Now, by triangular inequality, it holds

$$\begin{aligned} \|u - u_h\|_{\mathcal{E}} &\leq \|u - \tilde{u}\|_{\mathcal{E}} + \|\tilde{u} - u_h\|_{\mathcal{E}}, \\ &\leq \|\tilde{u} - u_h\|_{\mathcal{E}} + O(h), \text{ using (B.12) with } p = 2. \end{aligned}$$

So that

$$\left| \|u\|_{\mathcal{E}}^2 - \|u_h\|_{\mathcal{E}}^2 \right| \leq \|u - u_h\|_{\mathcal{E}}^2 + O(h)\|\tilde{u} - u_h\|_{\mathcal{E}} + O(h^2).$$

Using now (B.18), it yields

$$\left| \|u\|_{\mathcal{E}}^2 - \|u_h\|_{\mathcal{E}}^2 \right| \leq \|u - u_h\|_{\mathcal{E}}^2 + O(h^{1+p}),$$

which ends the proof, since $1 + p = 3/2$ for FEM and 2 for GFEM & SGFEM.

C Derivation of the stopping criteria for the iterative solvers

In this last annex, we derive the stopping criteria for the iterative solvers discussed in Section 6. The stopping criterion shown in (6.1) is derived as follows. First, let us assume, as indicated by a

priori error estimation, that there exists a constant $A > 0$, independent of the mesh size h , such that

$$h^p \leq A\epsilon^h,$$

where ϵ^h is the discretiation error, $p = 1/2$ for FEM and $p = 1$ for GFEM & SGFEM.

There exist [27] constants B, C , independent of the mesh size h , such that the following inverse estimates hold for all $v_h = \sum_{k \in \mathcal{N}_d^h} c_k N_k \in S^h$, denoting $\mathbf{c} = [c_k]_{k \in \mathcal{N}_d^h}$

$$\begin{aligned} \|v_h\|_{L^2(\Omega)} &\leq Bh\|\mathbf{c}\|_{l^2}, \\ \|v_h\|_{\mathcal{E}} &\leq Ch^{-1}\|v_h\|_{L^2(\Omega)}. \end{aligned}$$

So it holds

$$\|v_h\|_{\mathcal{E}} \leq BC\|\mathbf{c}\|_{l^2}.$$

Now, at some iteration i , we have an approximate solution $v_h^i = \sum_{k \in \mathcal{N}_d^h} c_k^i N_k$ to (3.6), and we can form the residual vector $\mathbf{r}^i = \mathbf{f} - \mathbf{A}_{11}\mathbf{c}^i$. Of course, the exact (discrete) solution $u_h = \sum_{k \in \mathcal{N}_d^h} c_k N_k$ solves the (discrete) residual equation and we have $\mathbf{r}^i = \mathbf{A}_{11}(\mathbf{c} - \mathbf{c}^i)$. Using the spectral radius of \mathbf{A}_{11}^{-1} , we have

$$\begin{aligned} \|\mathbf{c} - \mathbf{c}^i\|_{l^2} &= \|\mathbf{A}_{11}^{-1}\mathbf{r}^i\|_{l^2}, \\ &\leq \rho(\mathbf{A}_{11}^{-1})\|\mathbf{r}^i\|_{l^2}. \end{aligned}$$

Further, since the condition number of \mathbf{A}_{11} follows $\kappa_2(\mathbf{A}_{11}) = O(h^{-2})$ (we mention that the largest eigenvalue of \mathbf{A}_{11} is bounded independently of h), there exists a constant D , independent of the mesh size h , such that

$$\rho(\mathbf{A}_{11}^{-1}) \leq \frac{D}{h^2}.$$

As a result, if we compute $e^i = \frac{\|\mathbf{f} - \mathbf{A}_{11}\mathbf{c}^i\|_{l^2}}{h^2}$ and perform iterations until $e^i < \epsilon/k$, where $\epsilon = h^{1/2}$

(recall that this is the FEM case), we obtain

$$\begin{aligned}
\delta_i &= \|u_h - v_h^i\|_{\mathcal{E}}, \\
&\leq BC\|\mathbf{c} - \mathbf{c}^i\|_{l^2}, \\
&\leq BC\rho(\mathbf{A}_{11}^{-1})\|\mathbf{r}^i\|_{l^2}, \\
&\leq \frac{BCD}{h^2}\|\mathbf{f} - \mathbf{A}\mathbf{c}^i\|_{l^2}, \\
&\leq BCD e^i, \\
&< \frac{BCD}{k}\epsilon, \\
&< \frac{BCD}{k}h^{1/2}, \\
&< \frac{ABCD}{k}\epsilon^h.
\end{aligned}$$

Thus, iterations are performed until the truncation error δ_i is smaller than the discretization error ϵ^h , up to a proportionality factor controlled by k and independent of the mesh size h . The constants A, B, C and D are unknown and thus taking k large enough, δ_i could be made sufficiently smaller than the discretization error.

We further note that in practice, the “effective” condition number of \mathbf{A}_{11} is reduced from $O(h^{-2})$ to $O(h^{-1})$ thanks to the FMG preconditioner. As a result, it is sufficient to compute $e^i = \frac{\|\mathbf{f} - \mathbf{A}_{11}\mathbf{c}^i\|_{l^2}}{h}$ and perform iterations until $e^i < \epsilon/k$, where $\epsilon = h^{1/2}$.

The stopping criterion shown in (6.4) is derived as follows. First, since the outside iteration scheme yields a geometrical decrease of the truncation error $\delta_i = \|u_h - v_h^i\|_{\mathcal{E}}$, we make the following assumption: there exist positive constants $B_2 \geq B_1 \geq 0$ and $0 < q < 1$, all independent of i such that

$$B_1 q^i \leq \delta_i \leq B_2 q^i.$$

We further assume that these bounds are somewhat sharp and thus do not overlap from one iteration to the next, that is

$$B_2 q < B_1. \tag{C.1}$$

This is only required so that the truncation error effectively decreases at each iteration: $\delta_{i+1} < \delta_i$, which is the behavior observed in our numerical experiments.

Then, it follows by triangular inequality

$$\begin{aligned}
\|v_h^i - v_h^{i-1}\|_{\mathcal{E}} &\leq \|u_h - v_h^i\|_{\mathcal{E}} + \|u_h - v_h^{i-1}\|_{\mathcal{E}}, \\
&\leq B_2(1+q)q^{i-1},
\end{aligned}$$

and

$$\begin{aligned}\|v_h^i - v_h^{i-1}\|_{\mathcal{E}} &\geq \|u_h - v_h^{i-1}\|_{\mathcal{E}} - \|u_h - v_h^i\|_{\mathcal{E}}, \\ &\geq (B_1 - qB_2)q^{i-1}.\end{aligned}$$

Thus we have bounded $\|v_h^i - v_h^{i-1}\|_{\mathcal{E}}$ by above and below

$$C_1q^{i-1} \leq \|v_h^i - v_h^{i-1}\|_{\mathcal{E}} \leq C_2q^{i-1},$$

with $C_1 = B_1 - qB_2 > 0$ and $C_2 = B_2(1 + q) > 0$, also independent of i . The computable quantity $\|v_h^i - v_h^{i-1}\|_{\mathcal{E}}$ thus follows the same behavior with the number of iterations as the truncation error δ_i . Using the last three iterates we can estimate the common ratio and the scale factor of this geometric sequence. Similarly to (C.1), we will assume that there exist bounds $D_2 \geq D_1 \geq 0$ in (C.2) that are sharp and do not overlap from one iteration to the next, that is

$$\begin{aligned}D_1q^{i-1} &\leq \|v_h^i - v_h^{i-1}\|_{\mathcal{E}} \leq D_2q^{i-1}, \text{ and,} \\ D_2q &< D_1.\end{aligned}\tag{C.2}$$

Again, the second condition is only required so that at each iteration: $\|v_h^i - v_h^{i-1}\|_{\mathcal{E}} < \|v_h^{i-1} - v_h^{i-2}\|_{\mathcal{E}}$, which is also the behavior observed in our numerical experiments.

Let us now recall the stopping criterion e^i of (6.4)

$$e^i = \frac{1}{\frac{1}{\|v_h^i - v_h^{i-1}\|_{\mathcal{E}}} - \frac{1}{\|v_h^{i-1} - v_h^{i-2}\|_{\mathcal{E}}}}.$$

Now using the bounds in (C.2), it yields

$$\frac{1}{\frac{1}{D_1q} - \frac{1}{D_2}}q^{i-2} \leq e^i \leq \frac{1}{\frac{1}{D_2q} - \frac{1}{D_1}}q^{i-2}.$$

Iterations are performed until $e^i < \epsilon/k$, where $\epsilon = h$ (recall that these are the GFEM & SGFEM cases). When this happens, we have

$$\begin{aligned}\delta_i &\leq B_2q^i, \\ &\leq B_2\left(\frac{1}{D_1q} - \frac{1}{D_2}\right)q^2e^i, \\ &< \frac{B_2q^2}{k}\left(\frac{1}{D_1q} - \frac{1}{D_2}\right)\epsilon, \\ &< \frac{B_2q^2}{k}\left(\frac{1}{D_1q} - \frac{1}{D_2}\right)h, \\ &< \frac{AB_2q^2}{k}\left(\frac{1}{D_1q} - \frac{1}{D_2}\right)\epsilon^h.\end{aligned}$$

Thus, iterations are performed until the truncation error δ_i is smaller than the discretization error ϵ^h , up to a proportionality factor controlled by k and independent of the mesh size h . Again, the constants A, B_2, D_1, D_2 and q are unknown, and thus taking k large enough, δ_i could be made sufficiently smaller than the discretization error.

References

- [1] Y. Abdelaziz and A. Hamouine. A survey of the extended finite element. *Computers & Structures*, 86(11):1141–1151, 2008.
- [2] I. Babuška and U. Banerjee. Stable generalized finite element method (SGFEM). *Computer Methods in Applied Mechanics and Engineering*, 201–204:91–111, 2012.
- [3] I. Babuška, G. Caloz, and J. Osborn. Special finite element methods for a class of second order elliptic problems with rough coefficients. *SIAM J. Numer. Anal.*, 31(4):945–981, 1994.
- [4] I. Babuška, X. Huang, and R. Lipton. Machine computation using the exponentially convergent multiscale spectral generalized finite element method. *ESAIM: M2AN*, 48(2):493–515, 2014.
- [5] I. Babuška and R. Lipton. Optimal local approximation spaces for generalized finite element methods with application to multiscale problems. *Multiscale Model. Simul.*, 9(1):373–406, 2011.
- [6] I. Babuška and J. M. Melenk. The partition of unity finite element method. *Int. J. Numer. Meth. Engng.*, 40(4):727–758, 1997.
- [7] I. Babuška and J. Osborn. Generalized finite element methods: their performance and their relation to mixed methods. *SIAM. J. Numer. Anal.*, 20(3):510–536, 1983.
- [8] E. Béchet, H. Minnebo, N. Moës, and B. Burgardt. Improved implementation and robustness study of the X-FEM method for stress analysis around cracks. *Int. J. Numer. Meth. Engng.*, 64(8):1033–1056, 2005.
- [9] T. Belytschko and T. Black. Elastic crack growth in finite elements with minimal remeshing. *Int. J. Numer. Meth. Engng.*, 45(5):601–620, 1999.
- [10] T. Belytschko, R. Gracie, and G. Ventura. A review of extended/generalized finite element methods for material modeling. *Modelling and Simulations in Material Science and Engineering*, 17(4):43–74, 2009.
- [11] P. Bochev and R. B. Lehoucq. On the finite element solution of the pure Neumann problem. *SIAM Review*, 47(1):50–66, 2005.

- [12] W. L. Briggs, S. F. McCormick, and H. Van Emden. *A multigrid tutorial*. Siam, 2000.
- [13] Dassault Systèmes Simulia Corporation. *Abaqus, Version 6.11, Documentation*, 2011.
- [14] J. E. Dolbow. *An Extended Finite Element Method with Discontinuous Enrichment for Applied Mechanics*. PhD thesis, Northwestern University, 1999.
- [15] C. A. Duarte, I. Babuška, and J. T. Oden. Generalized finite element methods for three-dimensional structural mechanics problems. *Computers & Structures*, 77(2):215–232, 2000.
- [16] C. A. Duarte, O. N. Hamzeh, T. J. Liszka, and W. W. Tworzydło. A generalized finite element method for the simulation of three-dimensional dynamic crack propagation. *Computer Methods in Applied Mechanics and Engineering*, 190(15):2227–2262, 2001.
- [17] C. A. Duarte and J. T. Oden. An h-p adaptive method using clouds. *Comput. Methods Appl. Mech. Engrg.*, 139(1–4):237–262, 1996.
- [18] C. A. Duarte and J. T. Oden. H-p Clouds – An h-p meshless method. *Numer. Methods partial Differential Equations*, 12(6):673–705, 1996.
- [19] Y. Efendiev and T. Y. Hou. *Multiscale Finite Element Methods: theory and applications*. Springer, 2009.
- [20] T. -P. Fries. A corrected XFEM approximation without problems in blending elements. *International Journal for Numerical Methods in Engineering*, 75(5):503–532, 2008.
- [21] T. -P. Fries and T. Belytschko. The extended/generalized finite element method: An overview of the method and its applications. *Int. J. Numer. Meth. Engng.*, 84:253–304, 2010.
- [22] M. Griebel and M. A. Schweitzer. A particle-partition of unity method, Part VII: Adaptivity. In M. Griebel and M. A. Schweitzer, editors, *Meshfree Methods for Partial Differential Equations III*, volume 57 of *Lecture Notes in Computer Science and Engineering*. Springer, 2006.
- [23] V. Gupta, C. A. Duarte, I. Babuška, and U. Banerjee. A stable and optimally convergent Generalized FEM (SGFEM) for linear elastic fracture mechanics. *Comput. Methods Appl. Mech. Engrg.*, 266:23–39, 2013.
- [24] V. Gupta, C. A. Duarte, I. Babuška, and U. Banerjee. Stable GFEM (SGFEM): improved conditioning and accuracy of GFEM/XFEM for three dimensional fracture mechanics. *Comput. Methods Appl. Mech. Engrg.*, 289:355–386, 2015.
- [25] W. Hackbusch. *Multi-grid methods and applications*, volume 4. Springer-Verlag Berlin, 1985.

- [26] N. J. Higham. *Accuracy and Stability of Numerical Algorithms*. SIAM, Philadelphia, USA, 2002.
- [27] C. Johnson. *Numerical solution of partial differential equations by the finite element method*. Courier Corporation, 2012.
- [28] P. Laborde, J. Pommier, Y. Renard, and M. Salaün. High order extended finite element method for cracked domains. *Int. J. Numer. Meth. Engng*, 64(3):354–381, 2005.
- [29] Livermore Software Technology Corporation. *LS-DYNA, Users manual*, 2013.
- [30] S. Loehnert. A stabilization technique for the regularization of nearly singular extended finite element method. *Comput.Mech.*, 54(2):523–533, 2014.
- [31] J. M. Melenk. *On Generalized Finite Element Methods*. PhD thesis, University of Maryland, 1995.
- [32] J. M. Melenk and I. Babuška. The partition of unity finite element method: Basic theory and applications. *Comput. Methods Appl. Mech. Engrg.*, 139(1–4):289–314, 1996.
- [33] A. Menk and S. P. A. Bordas. A robust preconditioning technique for the extended finite element method. *Int. J. Numer. Meth. Engrg.*, 85(13):1609–1632, 2011.
- [34] N. Moës, M. Cloirec, P. Cartraud, and J. F. Remacle. A computational approach to handle complex microstructure geometries. *Comput. Methods Appl. Mech. Engrg.*, 192(28–30):3163–3177, 2003.
- [35] S. Nicaise, Y. Renard, and E. Chahine. Optimal convergence analysis for the extended finite element method. *Int. J. Numer. Meth. Engrg.*, 86(4–5):528–548, 2011.
- [36] J. T. Oden and C. A. Duarte. *Clouds, cracks and FEMs*. Citeseer, 1997.
- [37] H. Sauerland and T. P. Fries. The Stable XFEM for two phase flows. *Computers & Fluids*, 87:41–49, 2013.
- [38] M. A. Schweitzer. *A Parallel Multilevel Partition of Unity Method for Elliptic Partial Differential Equations*, volume 29 of *Lecture Notes in Computer Science and Engineering*. Springer, 2003.
- [39] M. A. Schweitzer. Stable enrichment and local preconditioning in the particle-partition of unity method. *Numer. Math*, 118(1):137–170, 2011.
- [40] F. L. Stazi, E. Budyn, J. Chessa, and T. Belytschko. An extended finite element with higher-order elements for curved cracks. *Computational Mechanics*, 31(1–2):38–48, 2003.

- [41] T. Strouboulis, I. Babuška, and K. Copps. The design and analysis of the generalized finite element method. *Comput. Methods Appl. Mech. Engrg.*, 181(1–3):43–69, 2000.
- [42] T. Strouboulis, K. Copps, and I. Babuška. The generalized finite element method: an example of its implementation and illustration of its performance. *Int. J. Numer. Meth. Engng.*, 47(8):1401–1417, 2000.
- [43] T. Strouboulis, K. Copps, and I. Babuška. The generalized finite element method. *Comput. Methods Appl. Mech. Engrg.*, 190(32–33):4081–4193, 2001.
- [44] N. Sukumar, D. L. Chopp, N. Moës, and T. Belytschko. Modeling holes and inclusions by level sets in the extended finite element method. *Comput. Methods Appl. Mech. Engrg.*, 190(46–47):6183–6200, 2001.
- [45] J. H. Wilkinson. *The Algebraic Eigenvalue Problem*. The Oxford University Press, 1988.
- [46] Q. Zhang, U. Banerjee, and I. Babuška. Higher order stable generalized finite element method. *Numer. Math.*, 128(1):1–29, 2014.